



Contents lists available at *Dergipark*

## Journal of Scientific Reports-A

journal homepage: <https://dergipark.org.tr/pub/jsr-a>



**E-ISSN: 2687-6167**

**Number 57, June 2024**

### **RESEARCH ARTICLE**

*Receive Date: 06.03.2024*

*Accepted Date: 23.06.2024*

# A hybrid approach to obesity level determination with decision tree and pelican optimization algorithm

Nagihan Yağmur<sup>a\*</sup>

*Kütahya Dumlupınar University, Faculty of Engineering, Department of Computer Engineering, Kutahya, Türkiye,  
ORCID:0000-0002-6407-4338*

## **Abstract**

Approximately 2 billion people in the world struggle with "obesity" and factors like eating lifestyle, habits, health conditions and mode of transport affect obesity. In this study, an artificial intelligence and machine learning-based model has been developed to predict obesity levels. It is proposed to create a hybrid model by combining the Decision Tree (DT) algorithm with the Pelican Optimization Algorithm (POA) on the obesity dataset of 2111 patients in SSggle. These models emphasize the critical role of parameters, aiming to achieve high performance. To solve the classification problem of multi-class obesity level determination, fuzzy logic-based parameter optimization is used to achieve high performance. While obesity rates are increasing worldwide, the study, which aims to globalize the parameters with the random discovery strategy of POA, is thought to be helpful for health professionals and decision-makers by successfully predicting obesity levels.

© 2023 DPU All rights reserved.

*Keywords:* Artificial intelligence, obesity, machine learning, decision tree, Pelican optimization algorithm, hybrid model

## **1. Introduction**

Within the medical realm, the discussion on categorizing obesity as a disease has concluded, establishing a unanimous consensus: Obesity has received official recognition as a medical condition [1], [2], [3]. This view has been supported by the World Health Organization (WHO) and more recently by the European Commission. Many

\* Corresponding author.

*E-mail address:* [nagihan.yagmur@dpu.edu.tr](mailto:nagihan.yagmur@dpu.edu.tr)

other health organizations have also endorsed this classification. A broad understanding of the scientific and health reasons behind defining obesity as a disease is now generally accepted in the medical community [1].

The surge in obesity rates has become a significant public health issue both nationally and globally [4], [5], [6]. Apart from imposing substantial financial burdens on the healthcare system, obesity leads to numerous detrimental health consequences, including cancers, endocrine disruptions, musculoskeletal issues, and a notable rise in premature deaths attributed to cardiovascular diseases [6], [7].

It is possible for anyone to gain weight and become obese in the context of his or her lifestyle, family history and a number of other factors according to the World Health Organization guidelines [5], [8]. The largest contributor to obesity, or the condition of uncontrolled weight gain, is excess fat [9]. There are several ways to measure excess fat in the human body [9]. Body fat is measured by a person's weight and height [9]. The most widely used scale is body mass index (BMI). Also used for measurements are skinfold thickness, bioimpedance, and waist circumference [9].

BMI furnishes crucial insights into categorizing body fat and body mass/weight, facilitating meaningful comparisons of weight across diverse groups. This helps to identify individuals at higher risk and prioritize interventions. Nevertheless, it is important to note that BMI fails to consider variations in body fat and muscle distribution, leading to the possibility that individuals with substantial muscle mass may register a high BMI even if their body fat percentage is low [9]. Waist size is another important measurement, especially used to assess overweight or obese individuals. However, it is important to understand that BMI alone is not sufficient in all cases and a more comprehensive assessment should be made by combining different measurements [9].

Obesity does not appear suddenly, but rather develops over time. It is important to apply machine learning techniques to assess a person's dietary habits and life choices. Data mining studies with the application of machine learning techniques will not only visualize changes in the follow-up of diet plans, but also allow this trend to predict negative events in the long term [9].

As engineering problems become more intricate, there is a rise in optimization challenges that classical optimization methods struggle to address. Therefore, new algorithms inspired by nature need to be developed. Among these artificial intelligence-based algorithms, there are methods inspired by nature such as POA, Butterfly Optimization Algorithm (BOA), Chicken Swarm Optimization Algorithm (CSO) [10], [11], [12].

Within data mining research, the exploration of influential attributes is commonly employed to evaluate classification outcomes. Feature selection aims to assess each parameter's impact while retaining those most critical to the model. Less impactful metrics can be excluded to focus on the most informative for categorization. Conventional logic-based methods restrict coefficients to binary values of zero or one. However, employing variably weighted coefficients permits a more nuanced assessment of each attribute's role. This finessed approach allows subtle adjustment of individual contributions to the predictive procedure. Such weighted evaluation proves vital for enhancing results and glean insights from large, complex datasets. Indeed, many published studies document the calculation of customized weights shed light on each factor's classification influence [13], [14], [15], [16], [17], [18].

The purpose of this research is to build a estimation model for assessing obesity levels based on machine learning techniques. This system accounts for an array of aspects including eating habits, lifestyle, medical conditions, and means of transportation. Here, we analyzed an obesity dataset from UCI comprising 2111 patients. The database categorizes individuals as Normal\_weight, Over\_weight\_1, Over\_weight\_2, Obesity\_type1, Obesity\_type2, Obesity\_type3 and Under\_weight. To assign obesity classes, we applied the 'fitctree' method utilizing the Decision Tree algorithm within Matlab [16]. Additionally, a hybrid approach is proposed combining DT with the POA. The POA outperforms other meta-heuristic algorithms in solving optimization problems by establishing a balanced relationship between exploration and exploitation. Derived from fuzzy logic principles, these blended models diligently appraise the importance of variables through the stochastic exploration of POA. In doing so, it strives to accomplish triumphant performance in categorizing individuals across multiple obesity levels.

The main contributions of this article can be summarized as follows:

- It comprehensively reviews the studies on obesity.
- It tests classical methods commonly used in the literature for obesity level classification or other classification problems with the same dataset.
- It proposes a new hybrid model to the literature by using DT and POA metaheuristics to classify the dataset to determine the obesity level.
- Inspired by the fuzzy logic approach, the proposed model better emphasizes the importance of parameters and achieves better results than the results obtained by classical methods.
- To the best of our knowledge, there is no other work on fuzzy logic-based parameter optimization to improve decision trees for multi-class obesity level classification.
- By improving the importance score calculated by the DT method for each parameter, an enhanced DT method is introduced to the literature.

## 2. Literature review

An investigation has recently been conducted into the use of 3D body scans to rectify deficiencies in obesity categorization through machine learning [19]. Duo-information (on body composition), consisting of Dual-energy X-ray absorptiometry, 3D body images, and Bioimpedance Analysis, was collected from one group of South Korean subjects. Following a machine learning framework based on 3D body scanning data for obesity classification, its utility was tested using a variety of metrics such as F1-score, Accuracy, Sensitivity and Precision. VKI and BIA are compared in this lightweight. VKI created Accuracy of 0.529, Sensitivity of 0.472, F1 score of 0.462, and Precision 0.458. On the other hand, BIA bettered these figures significantly with a Accuracy of 0.752, Sensitivity 0.742, Precision 0.751, and F1 score of 0.739. Not surprisingly, the proposed model compares favourably with both VKI and BIA, achieving an Accuracy of 0.800, a Sensitivity of 0.767, Precision 0.842 and F1 score 0.792. These values surpass the effectiveness of BMI and BIA in obesity categorization. In other a study, researchers endeavored to compare machine learning algorithms for a supervised task [20]. They employed Random Forest (RF) and K-Nearest Neighbor (KNN) models, submitting each to a battery of performance metrics. RF outpaced KNN following hyperparameter tuning, a 94% average accuracy. In a research investigation centered upon classifying obesity levels, dataset has been collected from the UCI Machine Learning Repository, which served as the basis for applying diverse machine learning methods. These included the Logistic Support Vector Machine (SVM), Regression, KNN, DT, and the sophisticated RF algorithm [21]. The analysis revealed the RF model delivered the most exceptional outcomes, attaining the highest measured accuracy at a %85.58. In a research initiative presenting a deep learning model designed to predict future obesity trends using the extensive medical records of children, a prominent pediatric health system in the United States provided an unmodified dataset [22]. Taking advantage of a versatile Long Short Term Memory (LSTM) network architecture, the model data in research consisted both static and dynamic Electronic Health Records (EHR) from 1-3 years prior. The model offers elaborate information to capture the onset of obesity among all persons between the ages of 3 and 20. Comparing the efficacy of the LSTM model with the data set from relevant literature studies, the model outperformed other models across most age groups-and this means it is very good at predicting patterns of obesity which depend on historical medical records. Similar findings led to another study that applied different machine learning techniques (including ensemble learning, generalized random forest, partial least square method, linear model) to the prediction of obesity trends based on age, height and weight, and BMI indicators [23]. The study achieved an impressive accuracy rate of more than 89%, showing that this comprehensive method predicts obesity patterns effectively.

When the studies using the same dataset are examined, the highest accuracy was obtained with J48 trees with 97.4% in the study using decision trees (J48), Bayesian networks (Naïve Bayes) and logistic regression [24]. In another study, chi-square, F-classification and mutual information classification algorithms were used to identify the most critical factors associated with obesity [25]. The performance of these models was compared using a neural network trained with different feature sets. Furthermore, the hyperparameters of the models are optimized with

Bayesian optimization techniques, which are faster and more efficient than traditional methods. The results showed that the neural network predicted the level of obesity with an average accuracy of 93.06%, 89.04%, 90.32% and 86.52% using all features and the features selected by the chi-square, F-classification and mutual information classification algorithms, respectively.

The processing of biomedical digital data and the evaluation of data records in hospitals are crucial for developing decision support systems for physicians [26]. Recently, numerous studies have focused on digital data processing and the classification of patient data records. In these studies, numerical data from patients, such as blood values, are typically processed to create systems that assist doctors in responding to new patients more quickly and accurately [27]. Both classical machine learning methods, such as Naïve Bayes, SVM, KNN, and regression, as well as deep learning techniques like Convolutional Neural Networks (CNN) and Stacked Autoencoders (SAE), have been utilized (knn nbayes cnn ekle). Additionally, nature-inspired metaheuristics, including Particle Swarm Optimization (PSO), Harris Hawks Algorithm (HHA), and Artificial Bee Colony (ABC), are increasingly being employed in these studies [13], [28], [29].

Feature selection Feature selection methods are important when guiding the impact of different parameters in a dataset on classification in data mining and machine learning, and to choose significant features when looking at data records [30]. It is characterized by getting rid of the less important parameters and keeping the more important ones in order to handle big datasets issues. We need to figure out 1) how to select features 2) and weight individual input parameters Page 129 So traditional feature selection someone receives a coefficient of 0 or 1 which means it is included accounted. But instead of selecting / discarding features we actually want to say that this feature is more significant compared to the other one by multiplying them by weight vector. This way, each parameter contributes to the success (the better) of classification to different degrees. However, it is better to give low weight rather than deleting a completely trivial parameter, with the creation of a function that is as insensitive to the hypothesis as possible (least dependence on this parameter while using its real range of values).

For many years, there have been numerous studies on the binary classification of obesity records (obese/not obese) and classification based on obesity levels. The literature proposes various machine learning methods for predicting and classifying obesity. These methods include KNN, SVM, RF, DT, and Learning Vector Quantization (LVQ), among others. These techniques may not always perform well across different datasets due to varying features and characteristics. Factors such as variability in the number of patient records significantly influence the success of obesity classification methods. Therefore, it is essential to develop new techniques to account for differences in dataset features and parameter numbers [31]. The increasing complexity of engineering problems has made solving them with traditional machine learning methods more challenging. The limitations of classical methods in addressing complex problems have necessitated the development of new nature-inspired techniques. Examples of nature-inspired algorithms include JAYA, CSO and Artificial Bee Colony (ABC). These methods are frequently used in data mining. Metaheuristics are successfully applied to various problems due to their ability to avoid local optima, derivative-free nature, flexibility, and ease of implementation. However, in recent years, the use of metaheuristics for obesity classification problems has been limited.

Additionally, although machine learning has been widely applied to obesity prediction and classification, there are few studies that comprehensively investigate various important factors such as an individual's health status and dietary habits and present a hybrid method. Simultaneously, emphasizing the significance of selected parameters to improve classification accuracy is a crucial factor in enhancing obesity prediction. Therefore, this paper introduces a hybrid model designed to identify classes including Normal\_weight, Over\_weight\_1, Over\_weight\_2, Obesity\_type1, Obesity\_type2, Obesity\_type3, and Under\_weight. The results of the proposed model are compared with those obtained from traditional machine learning techniques to provide a thorough comparative analysis.

Table 1. Dataset properties and types.

	<b>Parameter</b>	<b>Type</b>		<b>Parameter</b>	<b>Type</b>
1	Age	<b>Numerical</b>	20	CH2O_2	<b>1</b> : Between 1 liter and 2lt <b>0</b> : No
2	Gender	<b>0</b> : Male <b>1</b> : Female	21	CH2O_3	<b>1</b> : More than 2 liters <b>0</b> : No
3	Height (metre)	<b>Numerical</b>	22	SCC	<b>1</b> : Yes <b>0</b> : No
4	Weight (kg)	<b>Numerical</b>	23	FAF_1	<b>1</b> : Never <b>0</b> : Others
5	FHWO	<b>1</b> : Yes <b>0</b> : No	24	FAF_2	<b>1</b> : 1 or 2 times a week <b>0</b> : Others
6	FAVC	<b>1</b> : Yes <b>0</b> : No	25	FAF_3	<b>1</b> : 2 or 3 times a week <b>0</b> : Others
7	FCVC_1	<b>1</b> : Never <b>0</b> : Others	26	FAF_4	<b>1</b> : 4 or 5 times a week <b>0</b> : Others
8	FCVC_2	<b>1</b> : Sometimes <b>0</b> : Others	27	TUE_1	<b>1</b> : Less than 1 hour <b>0</b> : Others
9	FCVC_3	<b>1</b> : Always <b>0</b> : Others	28	TUE_2	<b>1</b> : 1 to 3 hours <b>0</b> : Others
10	NCP_1	<b>1</b> : 1 to 2 <b>0</b> : Others	29	TUE_3	<b>1</b> : More than 3 hours <b>0</b> : Others
11	NCP_2	<b>1</b> : 2 <b>0</b> : Others	30	CALC_1	<b>1</b> : No <b>0</b> : Others
12	NCP_3	<b>1</b> : More than 3 <b>0</b> : Others	31	CALC_2	<b>1</b> : Sometimes <b>0</b> : Others
13	NCP_4	<b>1</b> : No answer <b>0</b> : Others	32	CALC_3	<b>1</b> : Often <b>0</b> : Others
14	CAEC_1	<b>1</b> : No <b>0</b> : Others	33	CALC_4	<b>1</b> : Always <b>0</b> : Others
15	CAEC_2	<b>1</b> : Sometimes <b>0</b> : Others	34	MTRANS_1	<b>1</b> : Car <b>0</b> : Others
16	CAEC_3	<b>1</b> : Often <b>0</b> : Others	35	MTRANS_2	<b>1</b> : Motorcycle <b>0</b> : Others
17	CAEC_4	<b>1</b> : Always <b>0</b> : Others	36	MTRANS_3	<b>1</b> : Bicycle <b>0</b> : Others
18	Smoke	<b>1</b> : Yes <b>0</b> : No	37	MTRANS_4	<b>1</b> : Public transport <b>0</b> : Others
19	CH2O_1	<b>1</b> : Less than 1 liter <b>0</b> : No	38	MTRANS_5	<b>1</b> : Walking <b>0</b> : Others

When looking at the parameters in the Table 1, "FHWO" refers to "Family\_history\_with\_overweight," "FAVC" indicates "consumption of high-calorie food," "FCVC" represents "vegetable consumption frequency," "NCP" stands for "main meal frequency," "CAEC" signifies "food consumption between meals," "Smoke" pertains to "smoking information," "CH2O" denotes daily water consumption, "SCC" refers to "Information about whether or not you consume calories," "FAF" stands for "frequency of physical activity," "TUE" indicates "Time using technology devices," "CALC" represents "alcohol consumption," and "MTRANS" signifies "transportation."

### 3. Materials and methods

#### 3.1. Data preprocessing and dataset

For the investigation reported here, the source of the data set is UCI Laboratory, providing open access to its dataset [24], [32]. The dataset includes data from Mexico, Peru, and Colombia to estimate the obesity levels of individuals aged between 14 and 61 years. The data was collected using a web platform, through a questionnaire where anonymous users answered each question, and SMOTE (Synthetic Minority Over-Sampling Technique) was

used to balance the unbalanced dataset, resulting in a dataset of 17 attributes and 2111 records. However, since the variables are categorical, One Hot Encoding was applied to these variables and the categorical parameters were transformed into vector arrays consisting of 0 and 1. Thus, the number of parameters was increased to 38 as seen in Table 1. At the same time, due to the large differences between the values in the dataset, the data were normalized to 0-1 in order to make the data regular and comparable.

### 3.2. One hot encoding

By comparing each level of the categorical variable with a fixed reference level, it transforms a single variable with  $n$  observations and  $d$  different values into  $d$  binary variables, each with  $n$  observations [33]. The observations indicate the presence of the binary variable with 1 and its absence with 0 [33].

### 3.3. Decision trees

Decision trees (DTs) have gained widespread use in recent years for classification and pattern recognition, owing to their simplicity and straightforward rule-based construction [34]. DT adopts a multi-stage or sequential methodology in the classification process. A decision tree comprises three main components: nodes, branches, and leaves. Each attribute is represented by a node in the tree structure, with branches and leaves forming the additional components. The uppermost segment of the tree is denosssted as the root, while the lower portion is referred to as the leaf. The segments connecting the root to the leaves are termed branches [34], [35].

Building a decision tree entails asking a sequence of inquiries about the training data and drawing conclusions from the answers, with the objective of formulating decision rules. The classification process begins at the root node, progressing through nodes until branches or leaves without further subdivisions are encountered [34], [36]. To assess the generalization ability of the generated tree, a test dataset is employed. When a new test data point enters the tree structure established with the training data, it traverses through the tree, commencing from the root node. The new data, upon being tested at the root, is directed to a child node based on the test outcome, and this traversal continues until a specific leaf is reached. Each leaf corresponds to a single path or decision rule from the root [34], [36].

### 3.4. Pelican optimization algorithm

POA, conceptualized by Trojovský and Dehghani, draws inspiration from the foraging behavior of pelicans and stands as a metaheuristic optimization technique. Pelicans, renowned for their collective hunting approach, skillfully pinpoint their prey and execute a coordinated dive to capture it. These birds predominantly inhabit warm waters worldwide, favoring locations such as lakes, rivers, coasts, and marshes. Pelicans are social creatures, typically dwelling in groups, showcasing proficiency not only in flight but also in swimming [37]. Their keen eyesight and adept observation skills aid them during flight, and they primarily subsist on a diet of fish. When pelicans detect prey, they engage in a distinctive hunting behavior by running towards it from an altitude of 10-20 meters and directly diving into the water [38]. In the presence of fish schools, pelicans organize themselves into a line or U-shape formation, descending from the sky to manipulate the water with their wings. This strategic movement compels the fish to ascend, facilitating the pelicans in capturing them in their throat pouches. Pelicans' foraging behavior encompasses three distinct strategies [39].

The first strategy is the initialization step. Supposing that there are  $N$  pelicans in an  $M$ -dimensional space, the position of the  $i$ th pelican is  $P_i = [p_{i1}, p_{i2}, p_{i3}, \dots, p_{im}, \dots, p_{iM}]$ . Where  $p_{im}$  is  $i$ . Pelican's position in the  $m$ th dimension. Initially, the pelicans are randomly distributed in a certain range and the position update between  $up_m$  and  $low_m$  (0,1), which is the pelican's search range, is as follows.

$$P_{im} = low_m + random. (up_m - low_m) \quad i = 1, 2, \dots, M \quad (1)$$

The second strategy is to move towards its prey. In this phase, the pelican identifies its target and descends toward it from an elevated position. The adjustment of the pelican's location is formulated through Equation (2).

$$P_{im}^{t+1} = \begin{cases} P_{im}^t + rand. (S_m^t - \lambda \cdot P_m^t), & F(P_S) < F(P_i) \\ P_{im}^t + rand. (P_{im}^t - S_m^t), & F(P_S) \geq F(P_i) \end{cases} \quad (2)$$

Following the guidelines in Equation (2), where  $t$  signifies the current iteration count, the position of pelican  $i$  in dimension  $m$ , denoted as  $P_{im}^t$ , is expressed in relation to the position of the prey in dimension  $m$ , represented as  $S_m^t$ .  $F(P_S)$  designates the fitness function, and  $\lambda$  is a random value between 0 and 2, the third strategy step corresponds to the pelicans flapping their wings on the water's surface during hunting. This behavior is mathematically expressed in Equation (3).

$$P_{im}^{t+1} = P_{im}^t + \gamma \cdot \left(\frac{T-t}{T}\right) \cdot (2 \cdot random - 1) \cdot random \cdot P_{im}^t \quad (3)$$

In line with Equation 3, the existing iteration count denoted as  $t$ , alongside the maximum iteration count  $T$ , governs the  $\gamma$  neighborhood radius of  $P_{im}^t$ , expressed as  $\gamma \cdot \left(\frac{T-t}{T}\right)$ . This quantity is derived as a randomly generated value within the range of (0,1).

The steps of POA for solving the question of classification obesity levels are divided into as following:

1st step: The dataset was split into training and testing subsets, with  $xtrain$ ,  $ytrain$  for modeling and  $xtest$ ,  $ytest$  reserved for assessment.

2nd step: We initialized the pelican population and established iterations for optimization. The population size ( $N$ ) was set at 30 individuals and the algorithm run ( $IteN$ ) for 15 generations.

3rd step: The pelican population is initialized randomly by applying the formula in Equation 1. Since the pelican population will be the weight vector we are trying to optimize, its size should be ( $N \times 38$ ). Each row within the ( $N \times 38$ ) matrix with  $j = 1, 2, \dots, 38$  represents a weight vector, which aims to highlight the significance levels of values of parameters in the data, as stated in Equation 4.

$$d = [d_1, d_2, d_3, \dots, d_j] \quad (4)$$

4th step : As highlighted in Equation 5, the input data values are multiplied by the weight vector. Thus, significance of data parameter values is emphasized. In the study,  $i$  takes values between 1 and 2111.

$$Xtrain_{i,j} = xtrain_{i,j} * d_j \quad (5)$$

5th step: The MATLAB fitctree function is invoked, as outlined in Equation 6, to determine the fitness value linked to each member in the population, with a parameter setting of  $nTrees = 50$ .

$$Model = fitctree(XTrain, yTrain); \quad (6)$$

6th step : Classification is performed on the test data as in Equation 7.

$$yFound = Model.predict(xTest); \quad (7)$$

7th step: Make a comparison between  $y_{Found}$  and  $y_{test}$ . Then  $y$  is obtained by solving equation 8, which is an error value for such classifications. The hope here is to discover a weight vector leading to minimal  $y$ . Thus in these formulas Found represents the count of accurate records in the test dataset, and  $p_{Test}$  is a count of all records within the test dataset. The POA algorithm will be used to obtain the optimal solution for this weight vector.

$$y = 1 - \text{Found}/p_{Test} \tag{8}$$

8th step: A randomly chosen member from the population is designated as the prey.

9th step : When the fitness value of the prey falls below that of other individuals in the population, the individual's position undergoes an update and is subsequently stored as per Equation (2).

10th step: If the prey's value exceeds that of other individuals in the population, the individual's position is updated and recorded using Equation (3).

11th step: Evaluate the new location's fitness and compare it with the current fitness value. If the updated fitness is lower, integrate the current individual into the population, and adjust the pelican's position accordingly.

12th step: During the water surface flapping stage, recalculate and store the pelican's position based on Equation (4).

13th step: The fitness value of the new location value is compared with the old fitness value. If the new fitness value is smaller, the pelican location value is updated by adding the stored individual to the population.

14th step: Return to step 5 and repeat the process until the iteration concludes.

15th step: Ultimately, the optimal pelican location will be determined.

Table 2 outlines the values of variables crucial for identifying obesity classes.

Table 2. Parameter weights ( $d$ ) found by optimizing the DT method with POA.

Parameter	Weights	Parameter	Weights
Gender	0.430	CH2O_2	0
Age	0	CH2O_3	0.489
Height	0.228	SCC	0.286
WeighT	0.010	FAF_1	0.853
Family_history_with_overweight	0.330	FAF_2	0.190
FAVC	0	FAF_3	0
FCVC_1	0	FAF_4	0.081
FCVC_2	0	TUE_1	0.797
FCVC_3	0	TUE_2	0
NCP_1	0.134	TUE_3	0.402
NCP_2	0	CALC_1	0
NCP_3	0.640	CALC_2	0.459
NCP_4	0	CALC_3	0
CAEC_1	0	CALC_4	0
CAEC_2	0.693	MTRANS_1	0.248
CAEC_3	0.624	MTRANS_2	0
CAEC_4	0.264	MTRANS_3	0
Smoke	0	MTRANS_4	0.212
CH2O_1	0.098	MTRANS_5	0.270

Observing the weight vector associated with each parameter obtained through optimizing the decision tree method with POA in Table 2 reveals that certain parameters exert a predominant influence on classification success, whereas others exhibit values of 0 or 0-1.



#### 4. Evaluation

In this paper, a hybrid model is proposed by combining the POA metaheuristic with the decision tree method to detect Normal\_weight, Over\_weight\_1, Over\_weight\_2, Obesity\_type1, Obesity\_type2, Obesity\_type3 and Under\_weight classes. The results of the proposed model are compared with the results of classical machine learning methods.

The study implemented a 5-fold cross-validation method, where the dataset was partitioned into five subsets, and each subset was iteratively employed as a test set. Consequently, each subset served as a test set once, and all possible combinations were assessed, with the results subsequently averaged. For the evaluation of performance,

The purpose of ROC analysis is to evaluate the effectiveness of the results obtained from different methods and compare them using criteria such as Accuracy, Specificity, Sensitivity, F1-Score, Precision [18], [40], [41]. The main ROC parameters used in the analysis, FN (False Negative), TN (True Negative), FP (False Positive) and TP (True Positive), serve to indicate the accuracy of the classification results through true and false predictions. It is important to select appropriate metrics such as Accuracy, Precision, Recall and F1-Score to evaluate model performance. Accuracy measures the proportion of correct predictions, while precision shows how many of the values we predict as Positive are actually Positive. Recall is a metric that shows how many of the transactions we should have predicted as Positive we actually predicted as Positive, while F1-Score is the harmonic mean of Precision and Recall. Expressions for macro metrics, calculating the unweighted mean for each tag, are elucidated in Equations (9-12).

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (9)$$

$$\text{Recall}_m = \sum_c \frac{TP_c}{TP_c + FN_c}, c \in \text{snif} \quad (10)$$

$$\text{Precision}_m = \sum_c \frac{TP_c}{TP_c + FP_c}, c \in \text{snif} \quad (11)$$

$$\text{F1-Score}_m = \frac{2 \times \text{Precision}_m \times \text{Recall}_m}{\text{Precision}_m + \text{Recall}_m} \quad (12)$$

Within the equations, when considering the macro metric  $m$  and the classes = {Normal\_weight, Over\_weight\_1, Over\_weight\_2, Obesity\_type1, Obesity\_type2, Obesity\_type3 and Under\_weight classes}, the term  $TP_c$  denotes the count of samples accurately categorized as  $c$ .

#### 5. Experimental results

In this section, traditional methods frequently employed in the literature for obesity classification were subjected to tests using the "Classification Learner App" in Matlab [42]. The Linear Support Vector Machine, Quadratic Support Vector Machine, Cubic Support Vector Machine, Quadratic KNN models were assessed. As illustrated in Table 3, the Quadratic Support Vector Machine model demonstrated the highest accuracy rate, achieving an impressive 94.36%.

Table 3. Accuracy values of classical methods.

Model	Methods	Accuracy (%)
Model 1	LSVM	93.08
Model 2	QSVM	94.36



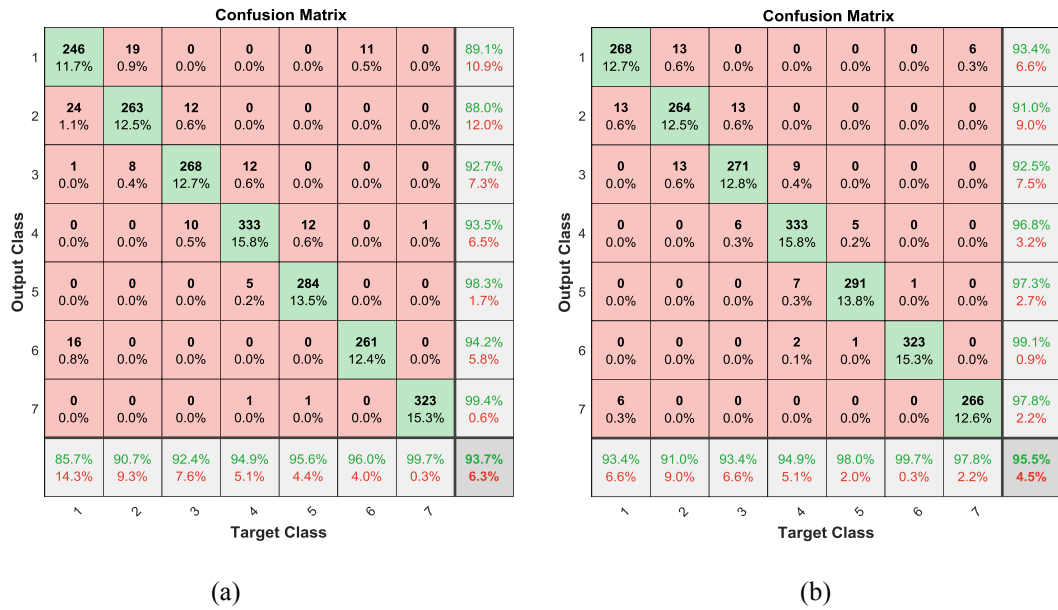


Fig. 2. (a) Confusion matrix of DT\_Model (b) Confusion matrix of POA\_DT\_Model.

Figure 1, where the classical methods are applied, shows that 145 patients were misclassified with LSVM, 118 with QSVM, 146 with QuSVM and 536 with QKNN. Similarly, when Table 3 is analyzed, it is seen that the accuracy achievements are QSVM, LSVM, QuSVM and QKNN, respectively.

To assess the efficacy of the proposed methods when solely applying the DT method, Figure 2 and Table 4 indicate that a collective total of 133 patients were misclassified, resulting in an accuracy success rate of 93.70%. With our proposed model, the accuracy was increased to 95.50% and only 95 patients were not classified correctly.

As a result, when the proposed model (POA\_DT\_Model) is compared with the model with only decision tree (DT\_Model), the accuracy is increased by 1.92% with the proposed model.

## 6. Results

The primary goal of this research was to craft decision support models for addressing the obesity epidemic ravaging our nation. In doing so, an assessment of crucial parameters including dietary habits, physical condition, physical activity levels, and lifestyle metrics is conducted to ascertain whether individuals are predisposed to obesity. Consequently, identifying individuals at elevated risk enables health professionals to formulate more targeted and effective intervention strategies. The utilization of obesity decision support systems facilitates the development of personalized health plans tailored to individuals' specific health conditions and requirements. This, in turn, contributes to a more effective realization of weight control and the pursuit of a healthy lifestyle for individuals.

This research combined the POA metaheuristic technique with decision tree modeling to shape a decision support system capable of categorizing obesity, which physicians can use to investigate clinical data. We believe that the focused hybrid strategies presented in this study effectively illustrate the importance in applying metaheuristics to optimize parameters as part of a successful computational optimization routine. Furthermore, we are convinced that our proposed hybrid framework based on fuzzy logic for parameter tuning would further enhance accuracy. Each

classification instance involved 5-fold cross validation. Success was appraised by examining a handful of performance metrics (e.g., accuracy, recall, precision, F-score). Results from our experiments showed the POA\_DT\_Model boasting 95.4% accuracy, providing an improvement of 1.92% over the DT\_Model. This relative success rate was 1.21% against classical methods. In conclusion, the hybrid approach proposed here is able to optimize very effectively the classification of obesity levels. We believe that the proposed models will improve their success performance, especially with increased parameter and data records, in future studies.

This study is suggested to be used in different data sets in the future. We believe that the proposed models will improve their performance in future studies, especially with increasing parameters and data records. The obtained results can be used to develop intelligent computational tools to determine the obesity levels of individuals and to create recommendation systems that monitor obesity levels. We believe that the hybrid approach logic presented in the literature will form the basis of different hybrid methods to be presented in the future, and parameter optimization with fuzzy logic approach will be a source of inspiration for different studies.

### Author Contribution

The writing of the manuscript and all analyses were performed by Nagihan Yağmur.

### Acknowledgements

In the course of this study, I would like to extend my gratitude to my cats, Karamel and Mahmutcan, for providing me with support and keeping my motivation high. Their love and loyalty strengthened me throughout this academic endeavor. Facing every challenge together made the time spent with them even more valuable.

### REFERENCES

- [1] M. Steele and F. M. Finucane, "Philosophically, is obesity really a disease?," *Obesity Reviews*, p. e13590, 2023.
- [2] T. K. Kyle, E. J. Dhurandhar, and D. B. Allison, "Regarding obesity as a disease: evolving policies and their implications," *Endocrinology and Metabolism Clinics*, vol. 45, no. 3, pp. 511–520, 2016.
- [3] A. M. Jastreboff, C. M. Kotz, S. Kahan, A. S. Kelly, and S. B. Heymsfield, "Obesity as a disease: the obesity society 2018 position statement," *Obesity*, vol. 27, no. 1, pp. 7–9, 2019.
- [4] CDC, "Overweight and Obesity." Accessed: Jan. 04, 2024. [Online]. Available: <http://www.cdc.gov/obesity/data/adult.html>
- [5] WHO, "Obesity and Overweight." Accessed: Jan. 04, 2024. [Online]. Available: <https://www.who.int/en/news-room/fact-sheets/detail/obesity-and-overweight>
- [6] V. Osadchiy *et al.*, "Machine learning model to predict obesity using gut metabolite and brain microstructure data," *Sci Rep*, vol. 13, no. 1, p. 5488, 2023.
- [7] J. J. Reilly and J. Kelly, "Long-term impact of overweight and obesity in childhood and adolescence on morbidity and premature mortality in adulthood: systematic review," *Int J Obes*, vol. 35, no. 7, pp. 891–898, 2011.
- [8] S. S. Shinde and R. S. Vaidya, "Automated Obesity Detection and Classification Via Live Camera Analysis" *International Research Journal of Modernization in Engineering Technology and Science*, vol. 5, no. 11, 2023.
- [9] S. A. Alsareii *et al.*, "Machine-Learning-Enabled Obesity Level Prediction Through Electronic Health Records," *Computer Systems Science and Engineering*, vol. 46, no. 3, pp. 3715–3728, 2023.
- [10] A. A. Heidari, S. Mirjalili, H. Faris, I. Aljarah, M. Mafarja, and H. Chen, "Harris hawks optimization: Algorithm and applications," *Future generation computer systems*, vol. 97, pp. 849–872, 2019.
- [11] X. Meng, Y. Liu, X. Gao, and H. Zhang, *A new bio-inspired algorithm: chicken swarm optimization*. Springer, p. 86–94.
- [12] A. Askarzadeh, "A novel metaheuristic method for solving constrained engineering optimization problems: crow search algorithm," *Comput Struct*, vol. 169, pp. 1–12, 2016.
- [13] N. Yagmur, I. Dag, and H. Temurtas, "A new computer- aided diagnostic method for classifying anaemia disease: Hybrid use of Tree Bagger and metaheuristics," *Expert Syst*, p. e13528, 2023.
- [14] N. Yagmur, I. Dag, and H. TEMURTAŞ, "A New Computer-Aided Diagnostic Method for Classifying Anemia Disease: Hybrid Use of Tree Bagger and Metaheuristics," *Authorea Preprints*, 2023.
- [15] S.-D. H. Ö. D. T. H. DÖRTERLER, "Hybridization of k-means and meta-heuristics algorithms for heart disease diagnosis," *New Trends in Engineering and Applied Natural Sciences*, p. 55, 2022.
- [16] S. Dörterler, H. Dumlu, D. Özdemir, and H. Temurtas, "Melezlenmiş K-means ve Diferansiyel Gelişim Algoritmaları ile Kalp Hastalığının

Teşhisi,” in *International Conference on Engineering and Applied Natural Sciences içinde (ss. 1840-1844)*. Konya, 2022.

- [17] S. Dörterler, “Kanser Hastalığı Teşhisinde Ölüm Oyunu Optimizasyon Algoritmasının Etkisi,” *Mühendislik Alanında Uluslararası Araştırmalar VIII*, p. 15, 2023.
- [18] S. Dörterler, H. Dumlu, D. Özdemir, and H. Temurtas, “Hybridization of Meta-heuristic Algorithms with K-Means for Clustering Analysis: Case of Medical Datasets,” *Gazi Mühendislik Bilimleri Dergisi*, vol. 10, no. 1, pp. 1–11.
- [19] S. Jeon, M. Kim, J. Yoon, S. Lee, and S. Youm, “Machine learning-based obesity classification considering 3D body scanner measurements,” *Sci Rep*, vol. 13, no. 1, p. 3299, 2023.
- [20] T. Turan, “Optimize Edilmiş Denetimli Öğrenme Algoritmaları ile Obezite Analizi ve Tahmini,” *Mehmet Akif Ersoy Üniversitesi Fen Bilimleri Enstitüsü Dergisi*, vol. 14, no. 2, pp. 301–312.
- [21] T. Cui, Y. Chen, J. Wang, H. Deng, and Y. Huang, “Estimation of Obesity levels based on Decision trees,” in *2021 International Symposium on Artificial Intelligence and its Application on Media (ISAIAM)*, IEEE, 2021, pp. 160–165.
- [22] M. Gupta, T.-L. T. Phan, H. T. Bunnell, and R. Beheshti, “Obesity Prediction with EHR Data: A deep learning approach with interpretable elements,” *ACM Transactions on Computing for Healthcare (HEALTH)*, vol. 3, no. 3, pp. 1–19, 2022.
- [23] K. Jindal, N. Baliyan, and P. S. Rana, “Obesity prediction using ensemble machine learning approaches,” in *Recent Findings in Intelligent Computing Techniques: Proceedings of the 5th ICACNI 2017, Volume 2*, Springer, 2018, pp. 355–362.
- [24] E. De-La-Hoz-Correa, F. Mendoza Palechor, A. De-La-Hoz-Manotas, R. Morales Ortega, and A. B. Sánchez Hernández, “Obesity level estimation software based on decision trees,” 2019.
- [25] F. H. Yagin *et al.*, “Estimation of obesity levels with a trained neural Network Approach optimized by the bayesian technique,” *Applied Sciences*, vol. 13, no. 6, p. 3875, 2023.
- [26] A. Clim, R. Zota, R. Constantinescu, and I. Ilie-Nemedi, “Health services in smart cities: Choosing the big data mining based decision support,” *Int J Healthc Manag*, vol. 13, no. 1, pp. 79–87, 2020.
- [27] E. Şahin, D. Özdemir, and H. Temurtas, “Multi-objective optimization of ViT architecture for efficient brain tumor classification,” *Biomed Signal Process Control*, vol. 91, p. 105938, 2024.
- [28] N. Yağmur, “Anemi Hastalığı Sınıflandırmasında Karga Arama Optimizasyon Algoritması,” in *Mühendislik Alanında Akademik Araştırma ve Derlemeler*, 2023, pp. 291–307.
- [29] N. Yağmur, I. Dag, and H. Temurtas, “Classification of anemia using Harris hawks optimization method and multivariate adaptive regression spline,” *Neural Comput Appl*, pp. 1–20, 2024.
- [30] R.-C. Chen, C. Dewi, S.-W. Huang, and R. E. Caraka, “Selecting critical features for data classification based on machine learning methods,” *J Big Data*, vol. 7, no. 1, p. 52, 2020.
- [31] S. Kilicarslan, M. Celik, and Ş. Sahin, “Hybrid models based on genetic algorithm and deep learning algorithms for nutritional Anemia disease classification,” *Biomed Signal Process Control*, vol. 63, p. 102231, 2021.
- [32] F. M. Palechor and A. de la Hoz Manotas, “Dataset for estimation of obesity levels based on eating habits and physical condition in individuals from Colombia, Peru and Mexico,” *Data Brief*, vol. 25, p. 104344, 2019.
- [33] K. Potdar, T. S. Pardawala, and C. D. Pai, “A comparative study of categorical variable encoding techniques for neural network classifiers,” *Int J Comput Appl*, vol. 175, no. 4, pp. 7–9, 2017.
- [34] T. Kavzoğlu and İ. Çölkesen, “Karar ağaçları ile uydu görüntülerinin sınıflandırılması,” *Harita Teknolojileri Elektronik Dergisi*, vol. 2, no. 1, pp. 36–45, 2010.
- [35] J. R. Quinlan, *C4. 5: programs for machine learning*. Elsevier, 2014.
- [36] M. Pal and P. M. Mather, “An assessment of the effectiveness of decision tree methods for land cover classification,” *Remote Sens Environ*, vol. 86, no. 4, pp. 554–565, 2003.
- [37] J. G. T. Anderson, “Foraging behavior of the American white pelican (*Pelecanus erythrorhynchos*) in western Nevada,” *Colonial Waterbirds*, pp. 166–172, 1991.
- [38] J. B. E. O’Malley and R. M. Evans, “Kleptoparasitism and associated foraging behaviors in American White Pelicans,” *Colonial Waterbirds*, pp. 126–129, 1983.
- [39] W. Tuerxun, C. Xu, M. Haderbieke, L. Guo, and Z. Cheng, “A wind turbine fault classification model using broad learning system optimized by improved pelican optimization algorithm,” *Machines*, vol. 10, no. 5, p. 407, 2022.
- [40] S. Kılıç, “Klinik karar vermede ROC analizi,” *Journal of Mood Disorders*, vol. 3, no. 3, pp. 135–140, 2013.
- [41] F. Aydemir and S. Arslan, “A System Design with Deep Learning and IoT to Ensure Education Continuity for Post-COVID,” *IEEE Transactions on Consumer Electronics*, 2023.
- [42] MATLAB, “Classification Learner.” Accessed: Jan. 04, 2024. [Online]. Available: <https://www.mathworks.com/help/stats/classificationlearner-app.html>