

# Entropy based C4.5-SHO algorithm with information gain optimization in data mining

G.Sekhar Reddy<sup>Corresp., 1</sup>, Ch. Suneetha<sup>2</sup>

<sup>1</sup> Dept of computer applications, Acharya Nagarjuna University, Guntur, Telangana, India

<sup>2</sup> Dept of computer applications, RVR&JC college of Engineering, Guntur, Telangana, India

Corresponding Author: G.Sekhar Reddy  
Email address: golamari.sekhar@gmail.com

Information efficiency is gaining more importance in the development as well as application sectors of information technology. Data mining is a computer-assisted process of massive data investigation that extracts meaningful information from the datasets. The mined information is used in decision-making to understand the behavior of each attribute. Therefore, a new classification algorithm is introduced in this paper to improve information management. The classical C4.5 decision tree approach is combined with Selfish Herd Optimization (SHO) algorithm to tune the gain of given datasets. The optimal weights for the information gain will be updated based on SHO. Further, the dataset is partitioned into two classes based on [b]quadratic entropy calculation and information gain. Decision tree gain optimization is the main aim of our proposed C4.5-SHO method. The robustness of the proposed method is evaluated on various datasets and compared with classifiers such as ID3 and CART. The accuracy and area under ROC (AUROC) parameters are estimated and compared with existing algorithms like ant colony optimization, particle swarm optimization and cuckoo search.

# Entropy Based C4.5-Sho Algorithm with Information Gain Optimization in Data Mining

G.Sekhar Reddy<sup>1</sup> and Ch. Suneetha<sup>2</sup>

<sup>1</sup>Research scholar, Acharya Nagarjuna University, Guntur

<sup>2</sup>Associate Professor, Dept of computer applications, RVR&JC college of Engineering, Guntur

Corresponding author:

G.Sekhar Reddy<sup>1</sup>

Email address: golamari.sekhar@gmail.com

## ABSTRACT

Information efficiency is gaining more importance in the development as well as application sectors of information technology. Data mining is a computer-assisted process of massive data investigation that extracts meaningful information from the datasets. The mined information is used in decision-making to understand the behavior of each attribute. Therefore, a new classification algorithm is introduced in this paper to improve information management. The classical C4.5 decision tree approach is combined with Selfish Herd Optimization (SHO) algorithm to tune the gain of given datasets. The optimal weights for the information gain will be updated based on SHO. Further, the dataset is partitioned into two classes based on quadratic entropy calculation and information gain. Decision tree gain optimization is the main aim of our proposed C4.5-SHO method. The robustness of the proposed method is evaluated on various datasets and compared with classifiers such as ID3 and CART. The accuracy and area under ROC (AUROC) parameters are estimated and compared with existing algorithms like ant colony optimization, particle swarm optimization and cuckoo search.

## INTRODUCTION

Information management is comprised of mining the information, managing data warehouses, visualizing the data and knowledge extraction from data and so on [Chen et al., 2018]. Consequently, different information management techniques are now being applied to manage the data to be analyzed. Hence, it is necessary to create repositories and consolidate data as well as warehouses. However, most of the data may be unstable; so it is essential to decide the data to be stored and discarded [Mohammad, Yin & Kasturi, 2019]. In addition, individual storage is required to manage real-time data to conduct research and predict trends. Data mining techniques are becoming more popular, recently getting attention towards rule mining methods such as link analysis, clustering and association rule mining [Asma et al., 2019]. Data mining discovers the substantial information, reasons and possible rules from huge datasets. It stands as an important source for information system based decision-making processes such as classification, machine learning and so on [Lin et al., 2019]. Data mining is generally a specific term to define certain computational analysis and results that comply with three main properties like comprehension, accuracy and user requirements. Data mining techniques are very useful while dealing with large datasets having vast amount of data. The data mining research community has been active for many years in analyzing various techniques and different applications of data mining [Swati, Hongmei & Karl, 2018].

A system that is combined with both data analysis and classification is suggested to create mining rules for several applications. For extracting the relevant information from systems, functional knowledge or rules automatically activates the mining process to provide rapid, real-time and significant operational basis. The classification approaches broadly used in data mining applications is efficient in processing large datasets [Xiaowei et al., 2018]. It maps an input data object into one of the pre-defined classes. Therefore, a classification model must be established for the given classification problem [Joao & Maria,

2019]. To perform the classification task, the dataset is converted into several target classes. The classification approach assigns a target type for each event of the data and allots the class label to a set of unclassified cases. This process is called supervised learning because all the training data are assigned as class tags. Therefore, classification is used to refer the data items as various pre-defined classes [Qi, Gengguo & Xiao, 2020]. The classifier is categorized into two approaches namely logical reasoning and statistical analysis. To create a well-trained classifier, training data is used to signify the key features of the classification problem under analysis [Xiangfei et al., 2020]. Once the classifier is trained, then the test dataset is evaluated by the classifier. The general performance of classification algorithm is relatively determined by the sensitivities to the minority target class. But in general, the predictions of the minority target class are optimal due to the primary design of the algorithms, which assume equal class distribution and application concepts [Solomon et al., 2019].

The most popular and simple classification technique is decision tree. As a general logical model, a decision tree repeats the given training data to create hierarchical classification [Fatima & Abdellatif, 2020]. It is a simplest form of classifier that can be stored densely and effectively in order to categorize the new data. It takes inputs in the form of training data set, attribute list and attribute selection method. A tree node is created by the algorithm in which attribute selection is applied to compute optimal splitting criteria. Then the final node generated is named based on the selected attributes [Irfan et al., 2019]. The training tuples subset is formed to split the attributes. Hence, parameter settings or domain knowledge is not needed for a decision tree. Moreover, it is capable of handling multidimensional data that offers good classification performance for common datasets [Phu et al., 2019]. Decision tree is also known as decision support tool which utilizes the model of tree-like or graph and the consequences are resource costs, utility and event outcomes [Jong, 2019]. In practical, the methods utilized to create decision trees normally produce trees with a low node factor and modest tests at each node. Also, the classifier contains different algorithms such as C4.5, ID3 and CART. The C4.5 algorithm is the successor of ID3, which uses gain ratio by splitting criterion for splitting the dataset. The information gain measure used as a split criterion in ID3 is biased to experiments with multiple outcomes as it desires to select attributes with higher number of values [Femando et al., 2019]. To overcome this, the C4.5 algorithm undergoes information gain normalization using split information value which in turn avoids over fitting errors as well.

In C4.5, the two criterions are used to rank the possible tests. The first criterion of information gain is to minimize the entropy of subsets and the second criterion of gain ratio is to divide the information gain with the help of test outcome information. As a result, the attributes might be nominal or numeric to determine the format of test outcomes. Decision trees are popular learning tool utilized in functional research, especially in results analysis to achieve a goal [Ludmila et al., 2019]. On the other hand, the C4.5 algorithm is also a prominent algorithm for data mining employed for various purposes. The generic decision tree method is created default for balanced datasets; so it can deal with imbalanced data too [Lakshmanaprabu et al., 2019]. The traditional methods for balanced data set when used for imbalanced datasets cause low sensitivity and bias to the majority classes. Some of the imbalance class problems include image annotations, anomaly detection, detecting oil spills from satellite images, spam filtering and software defect prediction, etc. [Fenglian et al., 2018]. The imbalanced dataset problem is seen as a classification problem where class priorities are very unequal and unbalanced. In this imbalance issue, a majority class has larger pre-probability than the minority class [Haoyue, MengChu & Qing, 2019]. When this problem occurs, the classification accuracy of the minority class might be disappointing [Xiaofen & Li, 2019]. Thus, the aim of the proposed work is to attain high accuracy in addition to high efficiency.

In data classification, accuracy is the main challenge of all applications. Information loss in dataset is problematic during attribute evaluation and so, the probability of attribute density is estimated. For this, the information theory called entropy based gain concept is utilized to enhance the classification task. Furthermore, common uncertainties of numerical data are used to measure the decision systems. A population based algorithm is utilized to optimize the gain attributes and to enhance the classification in complex datasets. The Selfish Herd Optimization (SHO) enhances the feature learning accuracy by effectively removing redundant features thereby providing good global search capability. The main contribution of the proposed work is summarized as follows.

- To solve the data classification problem using entropy based C4.5 decision tree approach and gain estimation.

- Selfish Herd Optimization (SHO) algorithm is utilized to optimize the information gain attributes of decision tree.
- The data is classified with high accuracy and AUROC of datasets is compared with existing techniques.

The organization of this paper is described as follows: introduction about the research paper is presented in Section 1, survey on existing methods and challenges are depicted in Section 2. The preliminaries are explained in Section 3. The working of proposed method is detailed in Section 4. Efficiency of optimization algorithm is evaluated in Section 5 and the conclusion of the proposed method is presented in Section 6.

## RELATED WORKS

Multiple learning process and multi-label datasets are widely used in different fields nowadays. [Anwar, 2019] evaluated the efficacy of Particle Swarm Classification (PSC) in data mining. PSC was utilized to design the classification model which classifies the queries into Bloom's taxonomy six cognitive-levels. Rocchio algorithm (RA) was used to mitigate the dimensionality of adverse effects in PSC. Finally, RA-based PSC was investigated with various feature selection methods for a scheme of queries. It is identified that the multi-label classification dealt with some problems where the classifier chain label order has a strong effect on the performance of classification. Nevertheless, it is too hard to find the proper order of chain sequences. Hence, [Lin et al., 2019] had proposed an ordering method based on the conditional entropy of labels where a single order was generated by this method. Reduced attributes can improve the accuracy of classification performances. The missed attribute values were typically not used in entropy or gain calculation. Information gain based algorithms tend to authenticate the attribute sets. Various measures were certainly affected from redundancy and non-monotonicity during attribute reduction. Therefore, a forward heuristic attribute reduction algorithm was proposed to solve the uncertainties in attribute selection. It simultaneously selects information attributes though unnecessary attributes were reduced in practice. [Can et al., 2019] proposed granular maximum decision entropy (GMDE) based on the measurement of monotonic uncertainty. Extreme decision entropy was developed in which the uncertainties of entropy are integrated with granulation knowledge. This investigation was validated with various UCI datasets and found to be computationally inexpensive.

The choice of dataset selection allows the extraction of highly representative information from the high-level data; so computational efforts were reduced among other tasks. A hybrid optimization based feature selection was proposed by [Rehab et al., 2019]. It is combined with slap swarm algorithm (SSA) and particle swarm optimization methods to enhance the efficacy of global and local search steps. his hybrid algorithm was examined on mixed datasets. It requires less time while the nodes quantity is reduced making it more desirable for large datasets. The SSA-PSO was employed to acquire best features from various UCI datasets. Also, redundant features were detached from the original datasets resulting in better accuracy. However, the accuracy is affected in complex datasets. To improve the classification performance of complex data, [Lin et al., 2019] introduced an attribute reduction method utilizing neighborhood entropy measures. The systems should have the ability to handle continuous data while maintaining its information on attribute classification. The concept of neighborhood entropy was explored to deal with uncertainty and noise of neighborhood systems. It fully reflects the decision-making ability by combining the degree of reliability with the coverage degree of neighborhood systems.

A clustering method based on functional value sequences has been proposed to accurately identify the functional equivalent programs with index variations. Because existing clustering programs were limited to structured metric vectors as in [Tiantian et al., 2020]. It is implemented for automated program repair to identify the sample programs from a large set of template programs. The average accuracy and average entropy were 0.95576 and 0.15497 respectively. However, the problem turned out to uncertain as the number of predictions is higher than the number of previous results. This issue was overcome by an alternative solution of priori weights and maximum entropy principle to attain the posteriori weights. [Carles et al., 2019] utilized a machine learning approach with single aggregated prediction from a set of individual predictions. A new factor presents a problem departing from the well-known maximal entropy hypothetical method and taking the distance among original and estimated integrated predictions. The suggested method was applied to estimate and measure predictive capability using prediction datasets.

It is difficult to perform feature selection (FS) for multi-label dimension curse in numerous learning

processes. Hence, [Mohsen et al., 2020] proposed a multi-label relevance–redundancy FS scheme based on Ant Colony Optimization (ACO) called ML-ACO. ML-ACO seeks to find the best features with lowest redundancy and many repetitions with class labels. To speed up the convergence, the cosine similarities between features as well as class labels are used as starting pheromone for each ant, and can be classified as a filter-based method. Various parametric entropies of decision tree algorithms are investigated by [Aldo, Juan & Luis, 2018]. Partial empirical evidences were provided to support the notion that parameter adjustment of different entropy activities influences the classification. Receiver operating characteristic (ROC) and Area under the ROC (AUROC) curve analysis provides an accurate criterion for evaluating decision trees based on parametric entropy. Various entropies such as Shannon entropy, Renyi entropy, Tsallis entropy, Abe entropy and Landsberg–Vedral entropy were discussed.

A new information classification algorithm has been introduced to improve the information management of restricted properties in [Hongbin et al., 2019]. Information management efficiency has gained more importance for the development of information technology through its expanded use. RLBOR algorithm was utilized to optimize the decision tree ratios. ID3 algorithm is a classical method of data mining that selects attributes with maximum information gain from the dataset at split node. However, decision tree algorithms have some drawbacks; it is not always optimal and it is biased in favor of properties that have higher values. In data classification, accuracy is the main challenge of all datasets. The resulting information loss is problematic for attribute evaluation while estimating the probability density of attributes. Due to the absence of classification information, it is challenging to perform potential classification. Consequently, an improved algorithm is utilized to solve the data classification issues.

## PRELIMINARIES

Entropy based measurements understands the decision system knowledge, its properties and some relations about the measurements. An optimization model is explored to enhance the performance of complex dataset classification. During prediction, the information gain optimal weights will be updated with the help of SHO algorithm. The nominal attributes of the dataset were designed by the ID3 algorithm. The attributes with missing values are not permitted. C4.5 algorithm, an extension of ID3 can handle datasets with unknown-values, numeric and nominal attributes [Gaurav & Hitesh, 2013]. C4.5 is one of the best learning based decision tree algorithms. In decision tree based classification, the training set is assumed as  $M$  and the number of training samples is mentioned as  $|M|$ . Here, the samples are divided into  $N$  for various kinds of  $K_1, K_2, \dots, K_n$  where the class sizes are labeled into  $|K_1|, |K_2|, \dots, |K_n|$ . A set of training sample is denoted as  $M$ , and the sample probability formula of class  $K_i$  is given in Equation (1).

$$p(M_i) = |K_i|/|M| \quad (1)$$

## Quadratic Entropy

Entropy is used to measure the uncertainty of a class using the probability of particular event or attribute. The gain is inversely proportional to entropy. The information gain is normally dependent on the facts of how much information was offered before knowing the attribute value and after knowing the attribute value. Different types of entropies are utilized in data classification. For a better performance, quadratic entropy is used in our work [Adewole & Udeh, 2018]. This entropy considers a random variable  $X$  as finite discrete with complete probability collection as mentioned in Equation (2).

$$p_i \geq 0 (i = 1, 2, \dots, n), \sum_{i=1}^k p_i = 1 \quad (2)$$

Here, the probability of event is denoted as  $p_i$ . The quadratic entropy of information is calculated by Equation (3).

$$Entropy M(x) = \sum_{i=1}^n p_i(1 - p_i) \quad (3)$$

Here,  $(M)$  specifies the information entropy of  $M$  (training sample set). For this particular attribute, the entropy of information is determined by Equation (4).

$$Entropy(M, H) = \sum_{g \in G} (|M_g|/|M|) * Entropy(M_g) \quad (4)$$

176 The entropy of attribute H is represented by Entropy (M, H), where H signifies attribute value. G denotes  
 177 all sets of values of g, and  $M_g$  denotes the subset of M, which is the value of H.  $|M_g|$  denotes the number  
 178 of elements in  $M_g$ , and number of elements of  $|M|$  in M.

### 179 Information Gain

The information gain is determined by Equation (5),

$$gain(M, H) = entropy(M) - entropy(M, H) \quad (5)$$

In a dataset M, Gain (M, H) denotes the information gain of attribute H. Entropy (M) signifies the sample set of information entropy and Entropy (M, H) denotes the information entropy of attribute H. In Equation (5), information gain is employed to find additional information that provides high information gain on classification. C4.5 algorithm chooses the attribute that has high gain in the dataset and use as the split node attribute. Based on the attribute value, the data subgroup is subdivided and the information gain of each subgroup is recalculated. The decision tree trained process is enormous and deep. Moreover, the gain ratio is different form information gain. It measures the information related to classification obtained on the basis of same partition. C4.5 uses the information gain and allows measuring a gain ratio. Gain ratio is described in Equation (6).

$$gain - ratio(M, H) = gain(M, H) / (split - info(M, H)) \quad (6)$$

Where,

$$split - info(M, H) = \sum_{g=1}^n -(M_g/M) \log_2(M_g/M) \quad (7)$$

180 The attribute with a maximum gain rate is selected for splitting the attributes. When the split information  
 181 tactics is 0, the ratio becomes volatile. A constraint is added to avoid this, whereby the information gain  
 182 of the test selected must be large at least as great as the average gain over all tests examined.

### 183 C4.5 decision tree

184 Ross quianlan developed the C4.5 algorithm to generate a decision tree. Many scholars have made various  
 185 improvements in the tree algorithm. However, the problem is that tree algorithms require multiple  
 186 scanning and deployment of data collection during the building process of decision trees. The processing  
 187 speed is slow because the dataset is too large to fit into the memory. The advantages of C4.5 algorithm are  
 considerable but a dramatic increase in demand for large data would improve its performance.

---

#### Algorithm 1: Pseudo code for C4.5 decision tree algorithm

---

```

Input: Dataset
Output: Decision tree
//Start
  for i=1: number of data
    Calculate information gain
  end
  HG= Attribute with highest information gain
  Tree = Create a decision node for splitting attribute HG
  New data= Sub datasets based on HG
  for i=1: number of New data
    Tree new= Create tree with new subset
    Attach tree to corresponding branch of Tree
  end
return
    
```

---

188 The C4.5 algorithm builds a decision tree by learning from a training set in which every sample is built on an attribute-value pair. The current attribute node is one and so the source node of the tree and the calculated extreme information gain ratio are gained in this manner [Suneetha & Raveendra, 2012]. It is

necessary to perform supervised data mining on the targeted dataset. This reduces the choice of classifiers in which a pre-defined classification could handle numerical data and classification. Each attribute is evaluated to find its ratio and rank during the learning phase of decision trees. Additionally, correlation coefficient is found to investigate the correlation between attributes as some dataset could not give any relevant result in data mining. In C4.5 decision tree algorithm, the gain is optimized by proposed SHO technique. The information gain is a rank based approach to compute the entropy. In this algorithm, the node with a highest normalized gain value is allowed to make decision, so there is a need to tune the gain parameter. The gain fitness is calculated based on the difference between actual gain value and new gain value. This is the objective function of the gain optimization technique which is described in Equation (8).

$$fitness = \min\{G - \hat{G}\} \quad (8)$$

Here,  $G$  and  $\hat{G}$  denotes actual and new gain respectively. Based on this fitness, the gain error is minimized by SHO and the gain value will be computed by using Equation (5). SHO can improve the learning accuracy, remove the redundant features and update the weight function of decision trees. The feature of SHO is random initialization generating strategy.

## PROPOSED METHOD: SELFISH HERD OPTIMIZATION (SHO)

SHO is utilized to minimize the gain error in a better way in optimization process. In meta-heuristic algorithms, SHO is a new branch inspired from group dynamics for gain optimization. It is instigated from the simulations of herd and predators searching their food or prey. The algorithm uses search agents moving in  $n$ -dimensional space to find solution for optimization problem. The populations of SHO are herd and predators where the individuals are known as search agents. In optimization areas, SHO is proved to be competitive with PSO for many tasks. The theory of Selfish Herd has been establishing the predation phase. Every herd hunts a possible prey to enhance the survival chance by accumulating with other conspecifics in ways that could increase their chances of surviving a predator attack without regard for how such behavior affects other individuals' chances of survival. This may increase the likelihood of a predator escaping from attacks regardless of how such activities disturb the survival probabilities of other individuals. The proposed SHO algorithm consists of different kinds of search agents like a flock of prey that lives in aggregation (mean of selfish herd), package of predators and predators within the said aggregate. This type of search agents is directed separately through fixed evolutionary operators which are centered on the relationship of the prey and the predator [Priyanka & Sankalap, 2020]. The mathematical model of SHO algorithm is given as follows.

### Initialization

The iterative process of SHO's first step is to initialize the random populations of animals as prey and predators thereby having one set of separable locations  $S = \{s1, s2, \dots, sN\}$ . Here, the population size is denoted by  $N$ . The position of animals is limited into lower and upper boundaries and the groups are classified into two, like prey and predator. Equation (9) is utilized to calculate the number of members in prey group.

$$n_p = \text{floor}(n \times \text{rand}(0.7, 0.9)) \quad (9)$$

Here, the quantity of prey group members is denoted as  $n_p$  where  $n$  denotes the population of the prey and the predators.

### Assignment of survival value

The survival value ( $SV$ ) of every animal is assigned and it is associated with the current best and worst positions of a known  $SV$  of whole population members. By optimization process, the present best and worst values are mentioned in the optimization problem. Then, the survival value will be determined by using Equation (10).

$$SV = f(x_i) - f_b / f_w - f_b \quad (10)$$

Where, worst and best fitness values are denoted by  $f_w$  and  $f_b$  respectively. Here,  $x_i$  represents the location of the prey or the predator.

## 215 Herd's leader movement

All herd members' movement is one of the significant steps in SHO. The location of leader of the herd is updated by Equation (11) as given in [Femando et al., 2017],

$$h_L = \begin{cases} h_L + 2 \times r \times \phi_{l.P_m} \times (P_m - h_m) & \text{if } SV_{h_L} = 1 \\ h_L + 2 \times r \times \psi_{l.y_{best}} \times (y_{best} - h_L) & \text{if } SV_{h_L} < 1 \end{cases} \quad (11)$$

216 Here, the tested selfish repulsion towards predators by current herd leader is denoted as,  $\phi_i$ , and  $r$  denotes  
217 the random number in the range (0, 1).  $\psi_L$  indicates the selfish attraction examined by the leader of the  
218 flock toward the global best location  $y_{best}$ .

Moreover, the location of the herd member  $h_a$  is updated based on two selections. Equation (12) is utilized to follow the herd and Equation (14) is utilized to recompense the group. Also, the selection is prepared based on some random variables.

$$h_a = h_a + f_a \quad (12)$$

Where,

$$f_a = \begin{cases} 2 \times (\beta \times \psi_{h_a, h_L} \times (h_L - h_a) + \gamma \times \psi_{h_a, h_b} (h_b - h_a)) & SV_{h_L} \leq SV_{h_u} \\ 2 \times \delta \times \psi_{h_i, h_m} \times (h_m - h_a) & \text{otherwise} \end{cases} \quad (13)$$

$$h_a = h_a + 2 \times \beta \psi_{h_L, y_{best}} \times (y_{best} - h_a) + \gamma \times (1 - SV_{h_a}) \times \hat{r} \quad (14)$$

219 Here,  $\psi_{h_a, h_m}$  and  $\psi_{h_a, h_L}$  indicates the selfish attractions examined through the herd member  $h_a$  towards  $h_b$   
220 and  $h_L$ , while  $\beta, \gamma$  and  $\delta$  indicates the random numbers in the range (0, 1). Also,  $\hat{r}$  represents the random  
221 direction unit vector.

## 222 Predator movement

The movement of every separable set of predators, the endurance of entities in the attacked flock and the distance between the predators from assault predators are taken into account in SHO. Based on the pursuit probability, the predator movement is determined as given in Equation (15).

$$P_i = \omega_{p_i, j_j} / \sum_{m=1}^{N_h} \omega_{p_i, j_j} \quad (15)$$

The prey attractiveness amongst  $p_i$  and  $h_j$  is denoted as  $\omega_{p_i, j_j}$ . Then the predator position  $X_p$  is updated by Equation (16).

$$X_p = X_p + 2 \times r \times (h_r - X_p) \quad (16)$$

223 In advance, each member of the predator and the prey group survival rate is recomputed by Equation (9).

## 224 Predation phase

The predation process is executed in this phase. Initially, the domain danger radius is computed by Equation (17).

$$R_r = \sum_{j=1}^n |y_j^l - y_j^u| \quad (17)$$

Where, upper and lower boundary members are represented by  $y_j^u$  and  $y_j^l$  respectively and the dimensions are denoted as  $n$ . After the radius calculation, a pack of targeted prey is computed by Equation (18).

$$T_{p_i} = h_j \in H | SV_{h_j} < SV_{p_i} | \|p_i - h_j\| \leq R_r, h_j \notin K \quad (18)$$

Here,  $SV_{h_j}$  and  $SV_{p_i}$  denotes the endurance tenets of  $P_i$  and  $h_j$  correspondingly.  $\|p_i - h_j\|$  signifies the Euclidean distance amongst the entities  $P_i$  and  $h_i$  respectively. The probabilities of the existence hunted are computed for every member of the set and is formulated in Equation (19) where  $K$  is  $\{K = K, h_j\}$ .

$$H_{p_i, h_j} = \omega_{p_i, h_j} / \sum_{(h_m \in T_{p_i})} \omega_{p_i, h_m}, h_j \in T_{p_i} \quad (19)$$



225 **Restoration phase**

Finally, the restoration is accomplished by making a set  $M = h_j \notin K$ . Here, K represents the set of herd member slayed for the duration of the predation phase. The mating probabilities are also determined by each member in Equation (20),

$$P_r = SV_{h_j} / \sum_{(h_m \in M)} SV_{h_m}, h_j \in M \quad (20)$$

226 Each  $h_j \in K$  is changed by a different result by SHO's mating operation, which is  $\text{mix}([h_{r1}, hr2, \dots, hrn])$ .  
 227 This SHO algorithm is utilized to optimize the gain function in data classification operation. Figure 1 displays the flow diagram of SHO algorithm.

---

**Algorithm 2:** Peseudo code for the proposed SHO algorithm in data classification

---

**Start**

Initialize the parametrs and locations of SHO by eq (9)

**for**

Each individual

Compute survival by eq (10)

**end for**

**while**  $K < K_{max}$

**for** every prey movement

**if** prey's leader

Update the location of prey leader by (11)

**else**

Update prey location by (12)

**end if**

**end for**

**for**

Every predator's movement

**for** each prey

Determine predation probability (13)

**end for**

Update predator location by (14)

**end for**

Re-compute survival value using eq (10)

Compute dangerous radius by (15)

Predation performance by (16) & (17)

Restoration performance by eqn (18)

K=K+1

**end while**

Global optimal output

Fitness for global optimal output

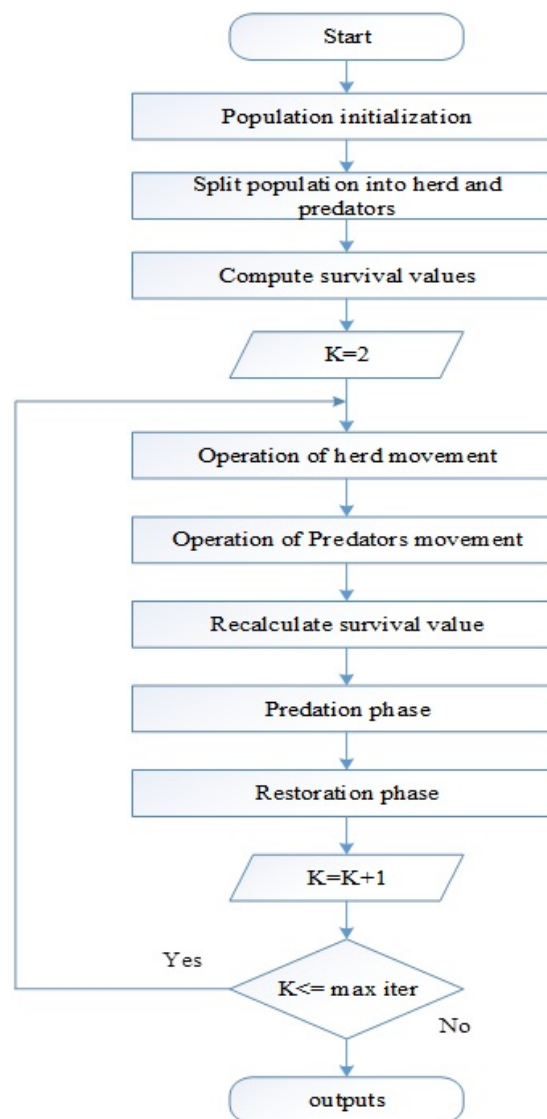
**end**

---

228

229 **RESULT AND DISCUSSION**

230 The efficiency of our proposed method is assessed by comparing its accuracy with other popular classifi-  
 231 cation methods like Particle Swarm Optimization (PSO) [Kun et al., 2014], Ant Colony Optimization  
 232 (ACO) [Femando, Alex & Colin, 2012], and Cuckoo Search (CS) Optimization. We estimated the  
 233 performance of proposed algorithm based on the accuracy as tested in 10 UCI datasets. The accuracy  
 234 of our proposed method is comparable to other optimization methods and various classifiers. We used  
 235 cross-validation in the datasets and got the mean value. The proposed method is greater than all existing  
 236 methods taken for comparison. SHO is combined with C4.5 classifier to produce greater accuracy than  
 237 a standard C4.5 classifier. The proposed decision tree classifier named C4.5-SHO is further compared  
 238 with C4.5, ID3 and CART. The description of ten data sets is tabulated in Table 1. These datasets



**Figure 1.** Flow diagram of SHO.

include Monks, Car, Chess, Breast-cancer, Hayes, Abalone, Wine, Ionosphere, Iris, and Scale [Aldo, Juan & Luis, 2018].

The proposed method is compared with existing entropies, optimization algorithms and different classifiers. The effectiveness is estimated based on the accuracy, AUROC and classifier.

# 1. Accuracy

The classification accuracy is measured based on Equation (21) [Kemal & Salih, 2009],

$$accuracy(A) = \sum_{i=1}^{|A|} assess(a_i) / |A|, a_i \in A \quad (21)$$

$$assess(a) = \begin{cases} 1, & \text{if } classify(a) = a.c \\ 0, & \text{otherwise} \end{cases} \quad (22)$$

Here, A is denoted as the data set to be classified (test set)  $a \in A$ ,  $a.c$  is the class of item  $a$  and  $classify(a)$  returns the classification through C4.5 classifier.

Data set	No of attributes	No of samples	Classes
Monks	7	432	2
Car	6	1728	4
Chess	6	28056	36
Breast-cancer	10	699	2
Hayes	5	160	3
Abalone	8	4177	2
Wine	13	178	2
Ionosphere	34	351	2
Iris	4	150	2
Scale	4	625	2

**Table 1.** Description of data set

In Table 2, the proposed C4.5-SHO decision tree classification accuracy is compared with other classifiers like C4.5, ID3 and CART. The accuracy of our proposed work is almost stable than the other. The accuracy of classification is depended on the training dataset. The accuracy of Iris data set is high (0.9986) compared to other data sets. The lowest accuracy of the proposed C4.5-SHO is 0.9437 in Scale data set. In comparison with existing classifiers, it is observed that the proposed approach has obtained a good accuracy.

Data set	C4.5-SHO	C4.5	ID3	CART
Monks	0.9832	0.966	0.951	0.954
Car	0.9725	0.923	0.9547	0.8415
Chess	0.9959	0.9944	0.9715	0.8954
Breast-cancer	0.9796	0.95	0.9621	0.9531
Hayes	0.9553	0.8094	0.9014	0.7452
Abalone	0.9667	0.9235	0.9111	0.9111
Wine	0.9769	0.963	0.9443	0.9145
Ionosphere	0.9899	0.9421	0.9364	0.9087
Iris	0.9986	0.9712	0.7543	0.8924
Scale	0.9437	0.7782	0.7932	0.7725

**Table 2.** Classification accuracy of the proposed classifier C4.5-SHO with C4.5, ID3 and CART

Data set	C4.5-SHO	C4.5-ACO	C4.5-PSO	C4.5-CS
Monks	0.9832	0.9600	0.9435	0.9563
Car	0.9725	0.9322	0.9298	0.9202
Chess	0.9959	0.9944	0.9944	0.9742
Breast-cancer	0.9796	0.9555	0.954	0.9621
Hayes	0.9553	0.90311	0.9322	0.9415
Abalone	0.9667	0.9500	0.9345	0.9247
Wine	0.9769	0.9240	0.8999	0.8924
Ionosphere	0.9899	0.9583	0.9645	0.9645
Iris	0.9986	0.9796	0.9741	0.9764
Scale	0.9437	0.9060	0.9177	0.8911

**Table 3.** Classification accuracy of proposed Algorithm with ALO, PSO and CS

In Table 3, the proposed C4.5-SHO decision tree classification accuracy is compared with other algorithms like ACO, PSO and CS. The accuracy of our proposed work is almost stable than the other. The accuracy of Iris data set is high (0.9986) compared to other data sets. The lowest accuracy of the proposed C4.5-SHO is 0.9437 in Scale data set. In comparison with existing algorithms, the proposed approach achieved good accuracy.

# 1. Area under ROC (AUROC)

The performance of classification model is shown through graph analysis of area under the Receiver Operating Characteristic curve (AUROC). This is dependent upon the attributes as well as classes. The

Data set	C4.5-SHO	C4.5	Id3	CART
Monks	0.9619	0.95713	0.9636	0.9791
Car	0.9819	0.9393	0.9891	0.8933
Chess	0.9673	0.9252	0.9090	0.9049
Breast-cancer	0.9793	0.9171	0.9730	0.9218
Hayes	0.9874	0.9069	0.9108	0.8360
Abalone	0.9647	0.9224	0.9573	0.9082
Wine	0.9914	0.9772	0.9497	0.9739
Ionosphere	0.9943	0.9680	0.9059	0.9560
Iris	0.9890	0.9048	0.7945	0.9481
Scale	0.9850	0.8562	0.7845	0.8007

**Table 4.** Area under the ROC curve of proposed C4.5-SHO with C4.5, Id3 and CART

Data set	C4.5-SHO	C4.5-ACO	C4.5-PSO	C4.5-CS
Monks	0.9935	0.9874	0.97668	0.9733
Car	0.98452	0.97908	0.97583	0.9659
Chess	0.99931	0.98612	0.9815	0.9503
Breast-cancer	0.9854	0.9795	0.9695	0.9581
Hayes	0.99616	0.92611	0.9442	0.9571
Abalone	0.9885	0.9828	0.9694	0.9566
Wine	0.9932	0.9830	0.8977	0.8964
Ionosphere	0.9954	0.9741	0.9630	0.9569
Iris	0.9873	0.9687	0.9656	0.9578
Scale	0.9858	0.9266	0.9165	0.8968

**Table 5.** Area under ROC curve of the proposed Algorithm with ALO, PSO and CS

proposed C4.5-SHO is compared with other classifiers like C4.5, ID3 and CART. The AUROC results are presented in Table 4 which shows that the proposed AUROC value is better than other algorithms.

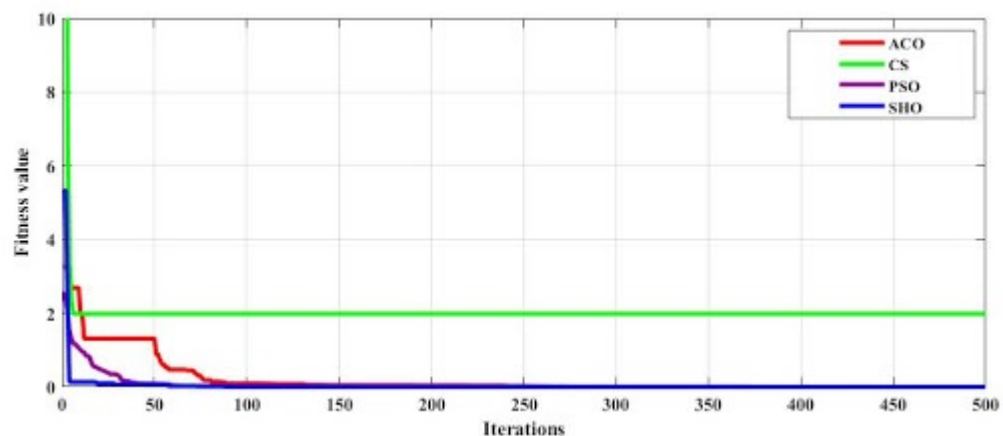
The proposed C4.5-SHO is compared with other optimization algorithms like ACO, PSO and CS. The AUROC results are presented in Table 5 which shows that the proposed AUROC value is better than existing algorithms. It is revealed that SHO not only reduces the complexity of decision trees but also enhances the accuracy.

# 1. Different entropy comparison

The proposed work follows Ray's quadratic entropy to get a better output. Based on the Ray's quadratic entropy, the information gain is optimized through the SHO algorithm. The entropy with SHO is compared to traditional SHO in terms of other entropies such as C4.5-SHO (Shanon entropy), C4.5-SHO (Havrda & charvt entropy), C4.5- SHO (Renyi entropy) and C4.5- SHO (Taneja entropy). Compared to other entropies, the Quadratic entropy achieved better accuracy in data classification for all data sets. Table 6 shows the entropy comparisons with proposed SHO.

Dataset	C4.5-SHO (Shanon entropy)	C4.5 – SHO(Havrda & charvt entropy)	C4.5 – SHO(Quadratic entropy)	C4.5- SHO (Renyi en- tropy)	C4.5- SHO(Taneja entropy)
Monks	0.9429	0.9756	0.9859	0.9926	0.9415
Car	0.9585	0.9527	0.9753	0.9895	0.9700
Chess	0.9510	0.9535	0.9907	0.9809	0.9401
Breast-cancer	0.9852	0.9558	0.9863	0.9564	0.9672
Hayes	0.9579	0.9460	0.9981	0.9476	0.9102
Abalone	0.9556	0.9618	0.9789	0.9715	0.9447
Wine	0.9485	0.9731	0.9823	0.9297	0.9317
Ionosphere	0.9319	0.9415	0.9665	0.9636	0.9036
Iris	0.9465	0.9807	0.9832	0.9514	0.9428
Scale	0.9725	0.8936	0.9747	0.9617	0.9031

**Table 6.** Entropy comparison

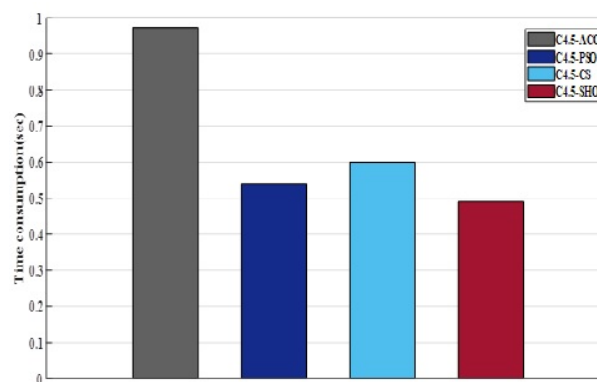


**Figure 2.** Comparison of convergence plot

Figure 2 illustrates the convergence plot of proposed SHO and similar existing algorithms. The proposed SHO achieved good convergence compared to existing techniques. The proposed work is based on gain optimization with SHO algorithm whereas the execution time is also the most important factor in data classification approach. On comparing the time-taken for analysis, the proposed method needs low computational time than the existing algorithms like ACO (0.974s), PSO (0.54s) and CS (0.6s). Table 7 and Figure 3 illustrate the computational time comparison.

Algorithm	Time(sec)
C4.5-ACO	0.974
C4.5-PSO	0.54
C4.5-CS	0.6
C4.5-SHO	0.49

**Table 7.** Computational Time



**Figure 3.** Comparison of computational time

278

## 279 CONCLUSION

280 Data mining is a broad area that integrates techniques from several fields including machine learning,  
 281 statistics, artificial intelligence, and database systems for the analysis of a large amount of data. This paper  
 282 presented a gain optimization technique termed as C4.5-SHO. The effectiveness of quadratic entropy  
 283 is estimated and discussed to evaluate the attributes in different datasets. This article presents the most  
 284 influential algorithms for classification. The gain of data classification information is optimized by the

proposed SHO algorithm. The evaluation of C4.5 decision tree based SHO results show that the AUROC is the best measure because of the classification of unbalanced data. The accuracy of proposed C4.5-SHO technique is higher than the existing techniques like C4.5, ID3 and CART. The proposed approach is compared with the algorithms of ACO, PSO and CS for AUROC. A better accuracy (average 0.9762), better AUROC (average 0.9909) and a better computational time (0.49s) are obtained from the gain optimized technique of C.5-SHO. In future, hybrid optimization technique is utilized to improve the data classification information gain.

## REFERENCES

- Adewole AP, Udeh SN. 2018.** The Quadratic Entropy Approach to Implement the Id3 Decision Tree Algorithm.
- Aldo RA, Juan BR, Luis MHS. 2018.** Statistical Entropy Measures in C4. 5 Trees. *International Journal of Data Warehousing and Mining (IJDWM)* 14:1-14.
- Anwar AY. 2019.** Swarm intelligence-based approach for educational data classification. *Journal of King Saud University-Computer and Information Sciences* 31:35-51.
- Asma E, Hasna N, Elkebir S, Ahmed H, Chafik N. 2019.** A novel information gain based approach for classification and dimensionality reduction of hyperspectral images. *Procedia computer science* 148:126-134.
- Can G, Zhihui L, Jie Z, Jiajun W, Wai KW. 2019.** Granular maximum decision entropy-based monotonic uncertainty measure for attribute reduction. *International Journal of Approximate Reasoning* 104:9-24.
- Carles B, Priscila E, Penélope H, Jose MP. 2019.** An entropy-based machine learning algorithm for combining macroeconomic forecasts. *Entropy* 21:1015.
- Fatima ES, Abdellatif H. 2020.** A MapReduce C4. 5 Decision Tree Algorithm Based on Fuzzy Rule-Based System. *Fuzzy Information and Engineering* 1-28.
- Fernando EBO, Alex AF, Colin GJ. 2012.** Inducing decision trees with an ant colony optimization algorithm. *Applied Soft Computing* 12:3615-3626.
- Fernando F, Erik C, Arturo V, Adrián G. 2017.** A global optimization algorithm inspired in the behavior of selfish herds. *Biosystems* 160:39-55.
- Fernando J, Carlos M, Enrico M, Jose TP, Gracia S, Guido S. 2019.** Multiobjective evolutionary feature selection for fuzzy classification. *IEEE Transactions on Fuzzy Systems* 27:1085-1099.
- Fenglian L, Xueying Z, Xiqian Z, Chunlei D, Yue X, Yu-Chu T. 2018.** Cost-sensitive and hybrid-attribute measure multi-decision tree over imbalanced data sets. *Information Sciences* 422:242-256.
- Gaurav LA, Hitesh G. 2013.** Optimization of C4. 5 decision tree algorithm for data mining application. *International Journal of Emerging Technology and Advanced Engineering* 3:341-345.
- Haoyue L, MengChu Z, Qing L. 2019.** An embedded feature selection method for imbalanced data classification. *IEEE/CAA Journal of Automatica Sinica* 6:703-715.
- Hongbin W, Tong W, Yucai Z, Lianke Z, Huafeng L. 2019.** Information classification algorithm based on decision tree optimization. *Cluster Computing* 22:7559-7568.
- Irfan SD, Agus PW, Anjar W, Sundari RA, Widodo S. 2019.** Decision Tree Optimization in C4. 5 Algorithm Using Genetic Algorithm. *In Journal of Physics: Conference Series* 1255: 012012.
- João RBJ, Maria DCN. 2019.** An iterative boosting-based ensemble for streaming data classification. *Information Fusion* 45:66-78.
- Jong-Seok L. 2019.** AUC4. 5: AUC-based C4. 5 decision tree algorithm for imbalanced data classification. *IEEE Access* 7:106034-106042.
- Kemal P, Salih G. 2009.** A novel hybrid intelligent method based on C4. 5 decision tree classifier and one-against-all approach for multi-class classification problems. *Expert Systems with Applications* 36:1587-1592.
- Kun-Huang C, Kung-Jeng W, Kung-Min W, Melani-Adrian A. 2014.** Applying particle swarm optimization-based decision tree classifier for cancer classification on gene expression data. *Applied Soft Computing* 24:773-780.
- Lakshmanaprabu SK., Shankar K, Ilayaraja M, Abdul WN, Vijayakumar V, Naveen C. 2019.** Random forest for big data classification in the internet of things using optimal features. *International*

*journal of machine learning and cybernetics* 10:2609-2618.

**Lin S, Xiaoyu Z, Jiucheng X, Shiguang Z. 2019.** An attribute reduction method using neighborhood entropy measures in neighborhood rough sets. *Entropy* 21:155.

**Lin S, Xiaoyu Z, Yuhua Q, Jiucheng X, Shiguang Z. 2019.** Feature selection using neighborhood entropy-based uncertainty measures for gene expression data classification. *Information Sciences* 502:18-41.

**Ludmila IK, Álvarez AG, José-Francisco D, and Iain ADG. 2019.** Instance selection improves geometric mean accuracy: a study on imbalanced data classification. *Progress in Artificial Intelligence* 8:215-228.

**Lin S, Xiao-Yu Z, Yu-Hua Q, Jiu-Cheng X, Shi-Guang Z, Yun T. 2019.** Joint neighborhood entropy-based gene selection method with fisher score for tumor classification. *Applied Intelligence* 49:1245-1259.

**Mohammad SA, Yin KC, Kasturi DV. 2019.** Identification of significant features and data mining techniques in predicting heart disease. *Telematics and Informatics* 36: 82-93.

**Mohsen P, Mohammad BD, Hossein N. 2020.** MLACO: A multi-label feature selection algorithm based on ant colony optimization. *Knowledge-Based Systems* 192:105285.

**Phu VN, Chau VTN, Tran VTHN, Dat ND. 2019.** A C4. 5 algorithm for english emotional classification. *Evolving Systems* 10:425-451.

**Priyanka A, Sankalap A. 2020.** A novel chaotic selfish herd optimizer for global optimization and feature selection. *Artificial Intelligence Review* 53:1441-1486.

**Qi X, Gengguo C, Xiao Z, Lei P. 2020.** Feature Selection Using Improved Forest Optimization Algorithm. *Information Technology and Control* 49: 289-301.

**Rehab AI, Ahmed AE, Diego O, Mohamed AE, Songfeng L. 2019.** Improved salp swarm algorithm based on particle swarm optimization for feature selection. *Journal of Ambient Intelligence and Humanized Computing* 10:3155-3169.

**Solomon HE, Mhd SS, Mamoun A, Ameer Al-N. 2019.** Variance ranking attributes selection techniques for binary classification problem in imbalance data. *IEEE Access* 7:24649-24666.

**Suneetha Ch, Raveendra BB. 2012.** Determining contribution of features in clustering multidimensional data using neural network. *I.J. Information Technology and Computer Science* 4: 29-36.

**Swati J, Hongmei H, Karl J. 2018.** Information gain directed genetic algorithm wrapper feature selection for credit rating. *Applied Soft Computing* 69:541-553.

**Tiantian W, KeChao W, Xiao HS, Lin L. 2020.** Data Mining in Programs: Clustering Programs Based on Structure Metrics and Execution Values. *International Journal of Data Warehousing and Mining (IJDWM)* 16:48-63.

**Wei C, Shuai Z, Renwei L, Himan S. 2018.** Performance evaluation of the GIS-based data mining techniques of best-first decision tree, random forest, and naïve Bayes tree for landslide susceptibility modeling. *Science of the total environment* 644:1006-1018. **Xiangfei M, Pei Z, Yan X, Hua X. 2020.** Construction of decision tree based on C4. 5 algorithm for online voltage stability assessment. *International Journal of Electrical Power & Energy Systems* 118:105793.

**Xiaofen T, Li C. 2019.** Artificial bee colony optimization-based weighted extreme learning machine for imbalanced data learning. *Cluster Computing* 22:6937-6952.

**Xiaowei G, Plamen PA, Ce Z, Peter MA. 2018.** A massively parallel deep rule-based ensemble classifier for remote sensing scenes. *IEEE Geoscience and Remote Sensing Letters* 15:345-349.