

Research Article

Semisupervised Bacterial Heuristic Feature Selection Algorithm for High-Dimensional Classification with Missing Labels

Hong Wang ¹, Yikun Ou ¹, Yixin Wang ¹, Tongtong Xing ¹ and Lijing Tan ²

¹College of Management, Shenzhen University, Shenzhen, China

²School of Management, Shenzhen Institute of Information Technology, Shenzhen, China

Correspondence should be addressed to Lijing Tan; mstlj@163.com

Received 28 September 2022; Revised 3 November 2022; Accepted 10 November 2022; Published 22 February 2023

Academic Editor: Fabio Caraffini

Copyright © 2023 Hong Wang et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Feature selection is a crucial method for discovering relevant features in high-dimensional data. However, most studies primarily focus on completely labeled data, ignoring the frequent occurrence of missing labels in real-world problems. To address high-dimensional and label-missing problems in data classification simultaneously, we proposed a semisupervised bacterial heuristic feature selection algorithm. To track the label-missing problem, a k -nearest neighbor semisupervised learning strategy is designed to reconstruct missing labels. In addition, the bacterial heuristic algorithm is improved using hierarchical population initialization, dynamic learning, and elite population evolution strategies to enhance the search capacity for various feature combinations. To verify the effectiveness of the proposed algorithm, three groups of comparison experiments based on eight datasets are employed, including two traditional feature selection methods, four bacterial heuristic feature selection algorithms, and two swarm-based heuristic feature selection algorithms. Experimental results demonstrate that the proposed algorithm has obvious advantages in terms of classification accuracy and selected feature numbers.

1. Introduction

The dimensionality of data, which consists of many features, is one of the most influential aspects of the classification model's effectiveness [1]. Based on feature properties, instances can be categorized into their respective classes. However, redundant, irrelevant, and noisy features in high-dimensional data will hamper classification accuracy [2, 3], e.g., medical or clinic data classification [4, 5]. Particularly, it is challenging to distinguish between representative and meaningless features without prior knowledge [6]. On the other hand, due to statistical norms and personal errors, data classification in real life often faces missing labels and loses more valid sample problems, which ultimately reduces the accuracy and robustness of the classification model [5, 7, 8]. To reduce feature dimensionality and improve the classification performance in classification, feature selection (FS) [9] is recommended to collect more relevant feature subsets

from the original data space. As a tool for optimizing data space, FS can make classification less complicated and improve the precision of classification models [10].

FS methods can be categorized as filter, wrapper, or embedded based on various feature evaluation criteria [11]. Filter methods use specific statistical metrics, such as information gain [12] and Fisher score [13], to evaluate the performance of created feature subsets, while wrapper methods use learning algorithms, such as K -nearest neighbor [14], naive Bayes [15], and linear discriminant analysis [16]. Embedded approaches, such as the least absolute shrinkage and selection operator [17] and ridge regression [18], embed FS into the training process for the learner. Filter methods typically execute faster than wrapper methods, but they cannot frequently achieve a higher degree of classification precision [19]. In addition, the process of designing embedded methods is intricate and necessitates plenty of prior experience [20]. Since the high-dimensional

classification problem with only partial labels is already a hard task, this research investigates wrapper-based FS methods to ensure a higher accuracy while avoiding increasing classification difficulties.

Wrappers seek to find the best subset from feature space according to one predetermined performance assessment. However, it is realistic to select all possible subsets of features measured by wrappers in high-dimensional classification problems because of the computational cost. Recently, wrappers based on population-based algorithms have been widely developed without the necessity of evaluating all possible subsets. Tran et al. [21] proposed the first variable-length particle swarm optimization representation for FS, enabling particles to have different and shorter lengths, which improves the performance of the algorithm. Considering the convergency of the population, Song et al. [22] proposed a variable-size cooperative coevolutionary strategy to optimize the searching population, which employs the idea of “divide and conquer” in the cooperative coevolutionary approach. However, since wrapper-based FS methods do not perform feature filtering in advance, the searching space for them is the whole data space [23]. This means that in high-dimensional classification tasks, their search space is very large, so they have to use a heuristic strategy like random search to reduce the cost of computation [24]. Nevertheless, classic heuristic strategy wrapper-based FS methods, such as particle swarm optimization-based FS [25], differential evolution-based FS [26], and genetic algorithm-based FS [27], do not account for all potential feature combinations [28].

In recent years, bacterial-based algorithms such as bacterial foraging optimization (BFO) [29] have been used to design FS methods to resolve combinatorial difficulties due to their global searching capability for control and optimization [30]. However, the intricate structure of BFO limits its computation efficiency. To achieve efficient classification results, bacterial colony optimization (BCO) [31] with a new bacterial life cycle was proposed and laid the foundation for the following research on bacterial-based FS algorithms and applications [6, 28, 30, 32, 33]. The majority of those research studies offer algorithmic enhancements in terms of weight setting, parameter optimization, and learning strategy optimization. Nonetheless, in actual applications, the integrity of data itself is a significant element, influencing the efficiency of FS, particularly the problem of incomplete sample labels, which is the most common and complicated task. This study focuses on developing an enhanced bacterial-based FS approach with a semisupervised learning strategy to address the high-dimensional medical data classification with partial labels.

According to the integrity of data labels, learning tasks can be categorized into supervised learning, semisupervised learning, and unsupervised learning [33]. In the supervised task, training data have complete label information, whereas in semisupervised learning, label information is only available in part. Unsupervised learning means that analyzed data do not contain labels [34]. In the absence of prior knowledge, the accuracy of supervised learning is generally higher than that of semisupervised learning. Nevertheless,

the cost of getting complete labeled data is extremely high in practical medical data collection. Moreover, unsupervised learning is usually used to disclose the initial pattern of unlabeled data [35]. Therefore, to address high-dimensional classification difficulty and label missing limitations in medical data, this paper investigated semisupervised medical data classification problems and optimized the classification model by learning from partially labeled data to classify unlabeled data into the correct class.

Semisupervised learning has been widely studied in different fields, and in the human activity recognition (HAR) problem, Chen et al. [36] designed a semisupervised deep learning model that is useful in solving the problem of imbalanced distribution of labeled data over classes from multimodal wearable sensory data. In video semantic recognition problems, Luo et al. [37] proposed a novel semisupervised feature selection method to learn the optimal graph, which aims to upgrade the performance of video semantic recognition. Since the research studies mentioned above are based on multimodal data, it makes more sense to employ deep learning or graph machine learning to overcome the problem of missing data labels. Despite the fact that these methods are effective for multimodal data, they incur substantial computational costs. Frequently, for a single mode of data, it is not necessary to use overly complex techniques. In contrast, feature selection methods based on wrappers need less computation, and hence, they are more suitable for single-type data. In terms of wrapper-based FS methods, certain representative semisupervised classification algorithms, such as ensemble SVM-based semisupervised FS [38] and rough set-based semisupervised FS [39], have demonstrated the ability of ensemble learning to solve label-missing problems. Nevertheless, these methods rely on ensemble classifiers to choose the best subset by voting for the results, which increases the computational cost marginally. As a straightforward and easy-to-use technique, K -nearest neighbor (KNN)-based semisupervised learning [40] offers great promise for improving the classification effect with missing labels. A number of studies utilize semisupervised KNN. Zhang et al. [38] demonstrated that the introduction of semisupervised learning with K -nearest neighbor (KNN) can enhance the available training sample size, provided that K is held constant. However, Mehta et al. [41] discovered that the magnitude of K would impact the efficiency of the algorithm. The precision of the results was enhanced by the use of an exhaustive procedure to determine a suitable value of K for solving the problem. Nevertheless, in partially labeled data, different learning densities of KNN may lead to biased results. When the value of K is small, model learning may not be comprehensive. When the value of K is large, operation cost may increase. In other words, the selection of the K value is a key issue to be explored. Therefore, this study attempts to develop a new semisupervised KNN learning approach that allows for the selection of K and can be combined with bacterial-based FS to form an effective classification method.

In this study, we propose a semisupervised bacterial heuristic feature selection (SHBFS) algorithm for the medical data classification mentioned earlier, i.e., label

incomplete and high-dimensional redundant features. The main contributions of this research are as follows:

- (i) A new self-adjusted semisupervised feature selection approach is proposed to solve the classification problems with missing labels and high-dimensional redundant features using a two-step self-training mechanism and an improved bacterial heuristic method
- (ii) The strategies of hierarchical population initialization, dynamic learning, and elite population evolution are proposed to enhance the capacity of the bacterial heuristic algorithm in searching for various feature combinations
- (iii) The proposed semisupervised bacterial heuristic feature selection algorithm is studied to be superior in addressing label incomplete and high-dimensional classification tasks in comparison to several state-of-the-art semisupervised FS algorithms

The rest of this paper is organized as follows: Section 2 gives the background of bacterial-based feature selection methods and some related works on these topics. The proposed method is introduced in Section 3. In Section 4, the experimental configuration is given. The experiments and analyses of the results are provided in Section 5. The final section presents the conclusions and a description of future work.

2. Related Work

The life cycle of the searching algorithm of the proposed bacterial-based FS approach in this study is inspired by BCO. Thus, this section briefly introduces its main principle and reviews of bacterial-based feature selection methods. More details are as follows.

2.1. Bacterial Colony Optimization. The life cycle of BFO is a triple-nested loop structure, which brings enormous computational complexity to solve high-dimensional problems. BCO simplifies the life cycle according to specific rules to address this computational drawback. Similar to BFO, BCO contains reproduction and elimination-dispersal processes. However, the chemotaxis steps in BCO are simplified as running and tumbling processes. Conditional controlling rules are used to cope with the traditional triple-nested loop structure to improve algorithm efficiency. The pseudocode of BCO is shown in Pseudocode 1.

2.1.1. Running Process. The running process is designed to speed convergence to the optimal position as

$$\theta_i^t = \theta_i^{t-1} + r_i \cdot (g \text{ best} - \theta_i^{t-1}) + (1 - r_i) \cdot (p \text{ best}_i - \theta_i^{t-1}), \quad (1)$$

where r_i shows the learning coefficient randomly generated between $[0, 1]$, $g \text{ best}$ is the best position in the current bacterial colony, and $p \text{ best}_i$ represents the individual optimal position during the chemotaxis process. In addition, communication schemes such as dynamic neighbor and group-oriented learning can be embedded into the running process.

2.1.2. Tumbling Process. The tumbling process avoids being trapped in the local optimum and explores more potential solution spaces. As shown in Equation (2), a random direction vector Δ_i is generated between $[-1, 1]$ for the i th bacterium. Although the randomly generated direction is correct, i.e., the current fitness is improved, the bacterium continues to exploit in the same direction.

$$\theta_i^t = \theta_i^{t-1} + r_i \cdot (g \text{ best} - \theta_i^{t-1}) + (1 - r_i) \cdot (p \text{ best}_i - \theta_i^{t-1}) + C(i) \cdot \frac{\Delta_i}{\sqrt{\Delta_i^T \cdot \Delta_i}}. \quad (2)$$

The reproduction and elimination mechanisms in BCO are consistent with those in BFO [29]. For the reproduction operation, half of the population with better performance is used to replace the remaining half with poor performance, while the elimination of BCO is realized by assigning the bacterium a new and random position within the search space. It can be formulated as follows.

If $P < \text{Ped}$, then

$$\theta^t = lp + \text{rand} \times (up - lp). \quad (3)$$

Otherwise,

$$\theta^t = \theta^t, \quad (4)$$

where Ped is a constant to determine the probability of the i th bacterium being assigned in a new position and up and lp are the upper and the lower boundaries of the search space, respectively. In this study, all BFO or BCO-based FS are referred to as bacterial-based FS, and the pertinent research reviews are detailed in the following section.

2.2. Bacterial Heuristic Feature Selection Methods. In recent years, research on bacterial-based FS algorithm improvements and applications has gradually increased. The process of bacterial-based FS consists primarily of the following steps [42]: (i) input the original data, (ii) randomly search for different features to form a small subset, (iii) use the subset to train the classifier, with the classification results to guide the bacteria search, (iv) output the optimal fitness of the

```

(1) Input: original data
(2) Initialization: P (population), MaxIt (max iterations), and C (chemotaxis step size)
(3) While maximum iterations are not satisfied do
(4)   For each bacterium do
(5)     Chemotaxis process (refer to Equation (1))
(6)     Fitness evaluation
(7)     If previous fitness < current fitness
(8)       Tumbling process (refer to (2))
(9)     End//alternative mechanisms
(10)    If the reproduction condition is satisfied do
(11)      Reproduction process (refer to article [29])
(12)    End
(13)    If elimination-dispersal conditions are satisfied do
(14)      Elimination-dispersal according to Equations (3) and (4)
(15)    End
(16)  End
(17) End//life cycle
(18) Output: optimal position with the best fitness

```

PSEUDOCODE 1: Bacterial colony optimization (BCO).

current iteration, and (v) loop steps ii~iv until the maximum number of iterations is reached and output the final optimal fitness.

In recent years, bacterial heuristic FS has many applications, including health care, recommendation, recognition, and model training [6, 30, 43, 44]. To improve the classification effect, bacterial-based FS has been improved in many ways. One improvement way is weight setting. Wang et al. [6] developed a weighted strategy to control the probability of different features being selected to enhance the accuracy. The other is population optimization, which can be further subdivided into position updates and population updates. For position updates, Wang and Chen [43] incorporated chaotic mechanisms into the chemotaxis and position-updating stages of bacterial populations to increase their adaptability. For population updates, some studies divided bacteria into multiple groups to perform different jobs under the control of different modified population updating strategies to improve the searching efficiency [32]. Furthermore, learning strategy optimization is also a common and useful method. For example, Kaur and Kadam [45] investigated multiobjective BFO to improve bacterial learning ability and improve the convergence speed of the algorithm. Wang et al. designed an adaptive attribute learning strategy to enhance the information communication ability among bacteria [30].

In summary, bacterial-based FS research focuses mostly on algorithm enhancement and the application of various situations. However, the combined effect of missing labels and high-dimensional redundant features poses significant challenges to optimizers (including bacterial heuristic algorithms) in FS, as the search space of FS problems expands exponentially and the proportion of incomplete data increases synchronously. Therefore, improving the effectiveness and efficiency of bacterial heuristic algorithms while considering semisupervised learning methods and data dimension reduction simultaneously is worth studying.

Therefore, in this study, we focus on the development of a semisupervised feature selection approach based on bacterial optimization to solve classification problems with missing labels and high-dimensional redundant features.

3. Proposed Approach

This section presents the proposed SHBFS approach to solve classification problems with missing labels and high-dimensional redundancy features. Figure 1 shows the structure of SHBFS. From the figure, we can see that the SHBFS approach consists of two main parts. On the one hand, a self-adjusted, semisupervised KNN strategy is presented for solving the problem of missing labels. On the other hand, an improved bacterial heuristic method for FS is presented for addressing the feature redundant problem, including three improvements: hierarchical population initialization, dynamic learning, and elite population evolution strategy. Hierarchical population initialization is used to obtain informative searching positions for bacteria to accelerate population convergence. Dynamic learning increases the searching variety of the algorithm by adaptively changing the search step length of bacteria. Finally, an elite population evolution strategy is employed to enhance the ability of bacteria to escape from the local optimum.

3.1. Self-Adjusted, Semisupervised KNN. The proposed self-adjusted, semisupervised KNN is a two-step self-training method, consisting of K value determination and label construction. Furthermore, to make the semisupervised learning method more adaptive to datasets of different sizes, the K value is adjusted as follows:

$$K = \begin{cases} 1: 1: NS, NS < 10, \\ 1: [\lg(NS)]: [10 * \lg(NS)], 10 \leq NS, \end{cases} \quad (5)$$

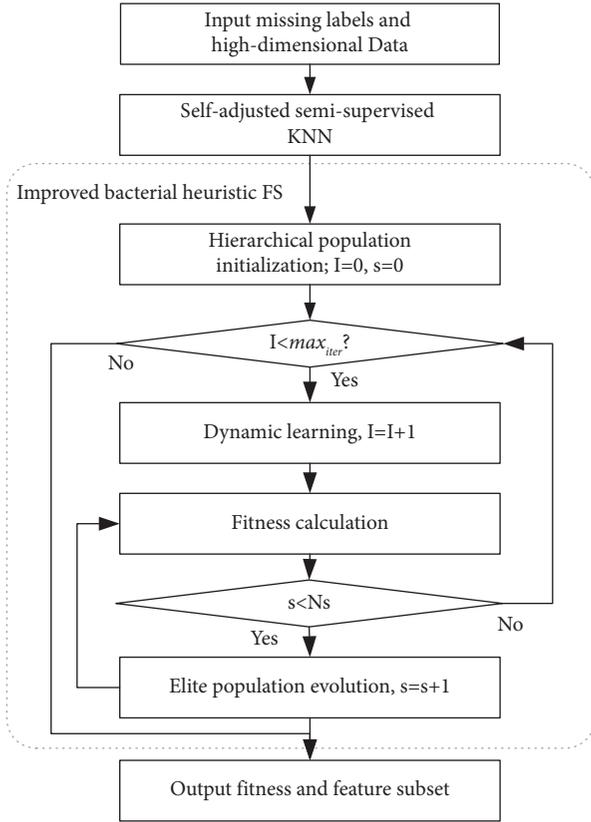


FIGURE 1: The overall structure of the proposed SHBFS approach.

where NS is the number of data samples, which means that the K value linearly increases when datasets have smaller samples, while the logarithmic function is employed for datasets with larger sample sizes. The primary process of the self-adjusted, semisupervised KNN is illustrated as follows:

Step 1. K Value Determination

The samples with the labeled class and unlabeled class are separately saved in the dataset L and dataset U . As mentioned previously, a self-adjusted, semisupervised KNN is presented for labeling the data with no assignment in categories and finding the best K value for classification. This step is to determine a K value for the label reconstruction using the labeled samples in L , provided in Pseudocode 2.

Step 2. Label Reconstruction

Label reconstruction is provided in Pseudocode 3. First, self-adjusted, semisupervised KNN is used to predict the labels for the samples from the dataset U . Then, newly labeled samples are moved from the dataset U to the dataset L . With increasing L in the space size, the self-training step can increase the learning efficiency of the training model.

3.2. Hierarchical Population Initialization. In BCO, the population is initialized randomly in a feasible space. However, addressing feature selection with high-dimensional features might make the bacterial colony fall into a poor searching position due to the high uncertainty of population initialization. As a result, more effort will be taken to jump out from their original position, which brings redundant computational complexity. To solve this problem, we develop a hierarchical population initialization strategy to enable bacteria to start at relatively good positions and further accelerate the convergence speed of the population. In contrast to the aforementioned variable-size cooperative coevolutionary technique, hierarchical population initialization does not use multipopulation for searching. Instead, it uses the idea of the proposed feature hierarchical division strategy to reconstruct a smaller search space before each search. The hierarchical population initialization consists of three steps. The details are as follows:

Step 1. Feature Ranking and Filtering

Initially, a symmetrical uncertainty (SU) [21] ranking is performed on the original features according to Equations (6)–(9). In this step, the correlations between features and classes are ranked, and the features' relevant significance is ordered from highest to lowest. After ranking, the worst 10 percent of features with significance below the mean are eliminated.

- (i) Symmetrical uncertainty (SU): the SU index has been widely used in traditional FS methods based on information theory. SU measures the uncertainty between feature variables $f \in F$, with label signals $l \in L$ given in Equations (6)–(9) based on the Shannon information entropy. In those formulas, $p(f)$ is the prior probability for all values of f and $p(f|l)$ is the posterior probability of f given l :

$$SU(f, l) = 2 \frac{H(f) - H(f|l)}{H(f) + H(l)}, \quad (6)$$

$$H(l) = - \sum_{i=1}^N p(l_i) \log_2 p(l_i), \quad (7)$$

$$H(f) = - \sum_{i=1}^N p(l_i) \log_2 p(l_i), \quad (8)$$

```

(1) Input: L (L is used to save the labeled samples in the original dataset)
(2) For each  $K$  obtained by Equation (5)
(3)   For each running time
(4)     Randomly divide  $L$  into two subsets/ * half with saved labels and half with removed labels */
(5)     Use the labeled subset to predict the labels of the unlabeled subset by KNN
(6)     Record the accuracy of label prediction on each running time
(7)   End
(8)   Record the  $K$  value and the average accuracy of label prediction of each  $K$ 
(9) End
(10) Output: Best ( $K$  value with the maximum average accuracy)

```

PSEUDOCODE 2: Determination of the K value.

```

(1) Input: labeled dataset L and unlabeled dataset U, best
(2) For each unlabeled sample in the dataset U
(3)   Use the samples from the dataset L to predict the label by KNN (using  $K_{best}$ )
(4)   Assign the predicted label to the unlabeled sample and move it from U to L
(5) End
(6) Output: the updated dataset L

```

PSEUDOCODE 3: Label reconstruction.

$$H(f|l) = - \sum_{j=1}^N p(l_j) \sum_{i=1}^N p(f_i|l_j) \log_2 p(f_i|l_j), \quad (9)$$

where $H(f)$ and $H(l)$ are the entropy of the feature variable f and the label signal l , respectively. N is the number of observation samples $x \in X$. $SU(f, l)$ evaluates the correlation between features f and label signals l . A larger $SU(f, l)$ indicates a higher significance of the feature f to the label l . This means that the feature f has more robust ability to discriminate labels, and the feature f needs to be selected into the feature subset.

Step 2. Feature Hierarchical Division

As shown in Figure 2, it is assumed that the numbers of SU are significant in the box. According to their significance, sorted features will be divided evenly into three layers, L_1 , L_2 , and L_3 . After that, 80% of the feature dimension will be randomly selected from the L_1 set, 15% from the L_2 set, and 5% from the L_3 set to form a searching position for bacteria. This strategy can exclude subpar features and shrink the search space when dealing with high-dimensional features.

Step 3. Feature Weight Updating

We assume that the feature size is H , and each feature of the i_{th} bacterium is denoted as f_i . We define the current fitness as $fit(f_i)$ and the historical fitness as $Fit(f_i)$. In this paper, we adopt a weight mechanism [6] to evaluate the performance of features. The rules are as follows: if $fit(f_i) < Fit(f_i)$, then the performance weight pf_i will be

increased by (12). Otherwise, pf_i will be decreased by Equation (13).

Given that, after completing the aforementioned procedure, there are still unselected features in each feature layer. To increase these features' probability of being selected in the future, we defined the unselected weight ($Uweight$) of f_i as uf_i and $Uweight = \{uf_1, uf_2, \dots, uf_H\}$. Then, the weight of each unselected feature will be updated by Equation (10) after Step 2. In each feature selection process, if one feature has been selected repeatedly in each search, then its uf_i will be decreased by Equation (11):

$$uf_i = \begin{cases} uf_i + 0.01uf_i, & uf_i \in L_1, \\ uf_i + 0.001uf_i, & uf_i \in L_2, \\ uf_i + 0.0001uf_i, & uf_i \in L_3, \end{cases} \quad (10)$$

$$uf_i = uf_i - (\max(uf_i) - \min(uf_i)), \quad (11)$$

$$pf_i = pf_i + \frac{|Fit(f_i) - fit(f_i)|}{Fit(f_i)}, \quad (12)$$

$$pf_i = pf_i - Fit(f_i) * |Fit(f_i) - fit(f_i)|, \quad (13)$$

where f_i is each feature of the i_{th} bacterium, uf_i is the unselected weight, pf_i is the performance weight,

Divide the features into three parts evenly
 D is sample; F is feature; L is the layer.

	F_1	F_2				F_H
D_1	0.98	0.88	0.51	0.58	0.26	0.11
D_2	0.94	0.83	0.45	0.66	0.28	0.17
	⋮		⋮			⋮
D_n	0.95	0.85	0.43	0.59	0.29	0.16
	└──────────┘		└──────────┘		└──────────┘	
	L_1		L_2		L_3	

FIGURE 2: Feature hierarchical division.

$\{L_1, L_2, L_3\}$ are layers obtained in Step 2, $fit(f_i)$ is the current fitness of f_i , and $Fit(f_i)$ is the historical fitness.

3.3. Dynamic Learning. In BCO, the chemotaxis step length of bacteria is governed by a set of fixed values denoted by $C(i)$. However, the lack of variation in step lengths may trap bacteria within the same search space. On a long-term basis, the diversity of feature subsets will decline. Therefore, in SHBFS, the running process is the same as that in BCO, while the tumbling process is improved by employing a dynamic learning strategy to increase the search variety.

Specifically, a dynamic learning strategy is proposed by adopting an adaptive chemotaxis step length changing strategy, which is denoted as aC [46], and a step length communication strategy dC . Equations (14) and (15) show that aC is affected by the bacterial size S , where $S = \{1, 2, \dots, i, i \in N^+\}$. We define the current fitness as fit , the upper bound of the step length as C^{ub} , and the lower bound of the step length as C^{lb} . ∂ is the disturbance factor. As the iteration proceeds, the disturbance effect of ∂ on aC will become small. In addition, the larger the fit , the larger the value of aC . The step length can be changed dynamically by aC :

$$\partial = \left| \left(1 - \frac{i}{S} \right) \times (C^{ub} - C^{lb}) + C^{lb} \right|, \quad (14)$$

$$aC = \frac{fit}{fit + \partial}. \quad (15)$$

There is no information communication among the bacteria in BCO. To enhance the convergence speed and improve the search capability, this paper presents a step length communication strategy. Let dC be the step length after communication, and its size is $S \times D$, where $D = \{1, 2, \dots, d, d \in N^+\}$ is the dimension of bacteria. Equation (16) shows the communication process of i_{th} bacteria in the t_{th} iteration:

$$dC^t = 0.01 \times aC + c_{pi} \cdot R_{pi} \cdot (p \text{ best}_i - \theta_i^t) + c_{gi} \cdot R_{gi} \cdot (g \text{ best}_i - \theta_i^t), \quad (16)$$

where θ_i^t is the current position of the bacterium, c_p and c_g are constant learning factors, and R_p and R_g are random disturbance terms. R_p and R_g are confined to $[0, 1]$. The step length size in SBHFS learns from the best population record of individual bacteria ($pbest$) and the best population record of the bacteria ($gbest$). For this, bacteria will prefer to learn from the record with a larger position excursion. After updating the step length dC , bacterial population tumbling is conducted as follows:

$$\theta_i^t = \theta_i^{t-1} + dC^{t-1} \cdot \frac{\Delta_i}{\sqrt{\Delta_i^T \cdot \Delta_i}} \cdot O_q, \quad (17)$$

where Δ_i is a random direction vector generated between $[-1, 1]$ for the i th bacterium. Due to varying data sizes, the range of bacterial location change is greater in large samples than in small samples. Therefore, $O_q = \{O_1, O_2, q = 1, 2\}$ has been proposed in this paper to adjust the offset of the bacterial position. In the tumbling process, $q = 1$; in the swimming process, $q = 2$. The setting of O_1 and O_2 is explained in Section 4.3. We define the number of features of the whole dataset as H , and the selected feature subset size is D . If $H < D$, $O_1 = O_2 = 1$.

After tumbling, the feature subset is formed by Equation (18), where $[\cdot]$ represents the rounding operator. The performance of the feature subset is measured by a classifier. Thus, we adopt the confusion matrix [47] as an evaluation metric, and the fitness is the error rate, which is updated by Equation (19).

$$\{[\theta_{i1}^t], \dots, [\theta_{id}^t]\}, \quad (18)$$

$$fit = \frac{FP + FN}{TP + TN + FP + FN}, \quad (19)$$

where FP is the false-positive result, FN is the false-negative result, TP is the true-positive result, and TN is the true-negative result. fit is the current fitness. Fit is defined as the historical fitness. The current best fitness is $fpbest$, and the historical best fitness is $fgbest$. The main process of dynamic learning is shown in Pseudocode 4.

3.4. Elite Population Evolution. In most bacterial-based methods, the population will randomly undergo dispersal elimination. This means that the new searching position of bacteria could be good or bad. The bad searching position may waste the search time. To make population evolution more meaningful, this paper designed an elite population evolution mechanism using $Pweight$ and $Uweight$ values aforementioned in Section 3.2, which are to guide bacteria to conduct reproduction and dispersal elimination.

In SBHFS, either reproduction or dispersal elimination will be conducted per iteration. The elite population evolution mechanism is proposed to determine which operation is executed, as depicted in Figure 3. After dynamic learning, bacteria will perform a swimming loop as BFO until they

```

(1) Input: fit,  $Pweight$ ,  $Uweight$ ,  $fpbest$ , and  $fgbest$ 
(2) For each bacterium
(3)   Running process by equation (1)//running
(4)   Update the guiding factor  $aC$  by Equations (14) and (15)
(5)   Update the step length  $dC$  by Equation (16)
(6)   Update the position by Equation (17)//tumbling
(7)   Get the feature subset by Equation (18)
(8)   Update fit by Equation (19) and update  $Uweight$  by Equation (11)
(9)   If fit <  $fpbest$ 
(10)      $fpbest = fit$ 
(11)     Update  $Pweight$  by Equation (12)
(12)   Else
(13)     Update  $Pweight$  by Equation (13)
(14)   End
(15)   If  $fpbest < fgbest$ 
(16)      $fgbest = fpbest$ 
(17)   End
(18)   Swimming by Pseudocode 5.
(19) End
(20) Output: fit and the feature subset

```

PSEUDOCODE 4: Dynamic learning.

```

(1) Input: fit, Fit,  $Pweight$ ,  $Uweight$ ,  $fpbest$ , and  $fgbest$ 
(2) For each bacterium
(3)   While doing the swimming loop
(4)     If fit < Fit
(5)       Fit =  $fi$ 
(6)       If  $bT >$  control threshold//it means bad effect searching
(7)         Do dispersal elimination
(8)         Rank  $F$  by  $Uweight$  and save as  $SF$ // $SF$  is the sorted feature
(9)          $\Theta =$  randomly selecting  $D$  features from the top half of  $SF$ 
(10)      Else// $D$  is the bacterial dimension
(11)        Do reproduction
(12)        For  $r = 1: Sr//Sr$  is the reproduction size of bacteria
(13)          Sort the position  $\Theta$  of bacteria
(14)          Sort features  $F$  by  $Pweight$  and save as  $SF$ 
(15)           $\Theta_{r+Sr} = SF(Sr)$ 
(16)        End
(17)      End
(18)    End
(19)    Get the feature subset by Equation (18)
(20)    Update fit by Equation (19) and update  $Uweight$  by Equation (11)
(21)    If fit <  $fpbest$ 
(22)       $fpbest = fit$ 
(23)      Update  $Pweight$  by Equation (12)
(24)      If  $fpbest < fgbest$ 
(25)         $fgbest = fpbest$ 
(26)      End
(27)    Else
(28)      Update  $Pweight$  by Equation (13)
(29)    End
(30)  End
(31) End
(32) Output: fit and the feature subset

```

PSEUDOCODE 5: Elite population evolution mechanism.

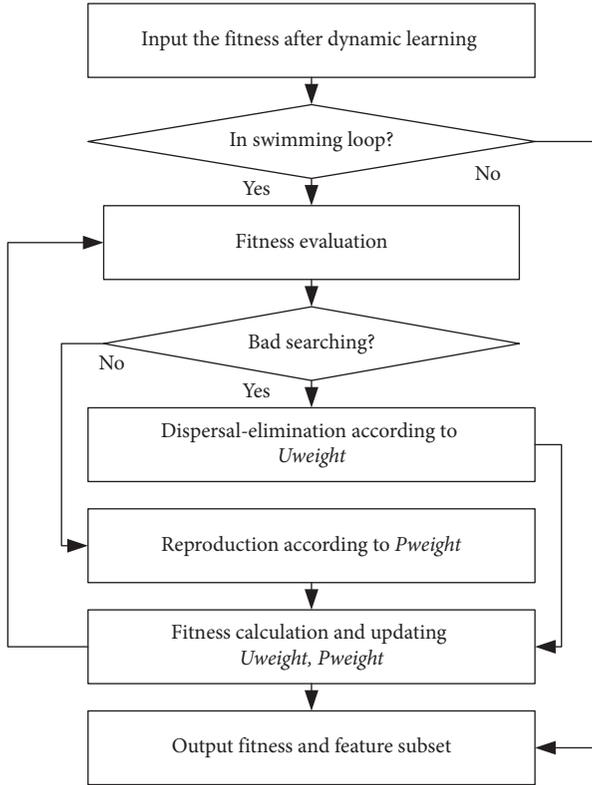


FIGURE 3: The elite population evolution mechanism.

meet the threshold N_s (see Table 1). In the swimming loop, each bacterium will first undergo a fitness evaluation to determine its performance. $errTre$ is defined as the performance threshold. If the fitness exceeds $errTre$, it will be counted in bT . When bT is larger than half of N_s , we can simply regard this bacterium by performing a bad search, and dispersal elimination is conducted based on the $Uweight$ matrix. Otherwise, bacteria will reproduce based on the $Pweight$ matrix. Next, we calculate the fitness of the new bacteria and updated two weights ($Pweight$ and $Uweight$). Finally, we repeat preceding steps until the end of the loop. The main process of the elite population evolution mechanism is given in Pseudocode 5.

The following is a description of reproduction and dispersal elimination processes: we assume that the bacterial dimension is D . In dispersal elimination, seldomly selected features are identified first by ranking features according to their $Uweight$. The new searching position for bacteria is then determined by randomly selecting D features from the top half of the seldomly selected features. In reproduction, features are initially ranked by their $Pweight$ to identify the highest-performing features of the search history. Then, the half population of bacteria with poor performance will be gradually replaced by the dimensions of the highest-performing features. The overall pseudocode of SBHFS is given in Pseudocode 6, and here, we analyze the computation time of feature selection (lines 3 to 13) of SBHFS. We suppose that there are S bacteria in the population, the max iteration time is I , the original number of features is D , and the swimming time is M and $M \ll I$.

First, we analyze the initialization part. The feature ranking by SU and feature weighting are the main time consumption step. The time complexity of calculating SU scores and weight for features are both $O(D)$, which are related to the number of features. Thus, the time complexity of the initialization part is $O(D) + O(D) \cong O(D)$. We further analyze the time computation of the main loop of Pseudocode 6 from lines 4 to 13. At iteration I , if the dynamic learning step in line 6 is conducted, the time complexity of this step is $O(SN_I)$ according to Pseudocode 4, where N_I is the selected features at iteration I and $N_I \leq D$. In the elite population evolution part in line 11 of Pseudocode 6, if the evolution choice is dispersal elimination, the time complexity is $O(SM)$. If the evolution choice is reproduction, the time complexity is $O(SM) + O(s) \cong O(SM)$, where s is the population to be updated and $s \leq S$ (usually, s is half of the population). Therefore, in the worst case for one iteration I , the complexity of SBHFS is $O(D) + O(SN_I) + O(SM) \cong O(SN_I)$. Since N_I denotes the number of selected features at iteration I and N_I is smaller than or equal to D , the time complexity of SBHFS at iteration I is smaller than or equal to $O(SD)$. Thus, we can reach the conclusion that the time complexity of the main loop of SBHFS during I iterations is not larger than $O(ISD)$.

4. Experimental Configuration

In this section, detailed information on the datasets, benchmark methods, and experimental design is given.

4.1. Datasets. In this paper, we verified the proposed method on different datasets, consisting of five high-dimensional microarray datasets and three benchmark datasets [6]. The description of the selected datasets is given in Table 2. #Features define the number of original features, #instance denotes the number of samples, and the number of classes is given in #Class. #Smallest class is the size of the class with the fewest instances, whereas #Largest class is the size of the class with the most instances. Among these datasets, Colon, SRBCT, DLBCL, Leukemia-AllAML (LA), and Central Nervous System (CNS) are datasets with the highest number of features up to 7129. All feature values in those five datasets are normalized within $[0, 1]$. Besides, the number of instances relative to the number of features in the last five datasets is considerably lower. Furthermore, all datasets are significantly imbalanced. These traits present FS and classification with formidable challenges. Since the proposed method is intended to handle missing label data, the original data will be transformed into partially labeled data, as described in Section 4.3.

4.2. Comparison Methods. The proposed SBHFS method is measured and compared with six recently widely recognized bioinspired wrapper FS algorithms, denoted as benchmark methods. The parameters of the comparison algorithms are shown in Table 1.

The adaptive chemotaxis bacterial foraging optimization algorithm (ACBFO) [42], improved swarming and

TABLE 1: The settings of parameters for benchmark methods.

Algorithms	Parameter settings
SBHFS	$errTre = 0.6, c_p = 0.0015, c_g = 0.0015, N_s = 4$ $C_{ub} = 0.15, C_{lb} = 0.05, O_1 = 1000, O_2 = 100$
ACBFO	$Nre = 5, Ned = 2, Nc = 10, Ns = 4, \alpha = 0.2,$ $Ped = 0.25, d_{attract} = h_{repellant} = 0.1, w_{attract} = w_{repellant} = 0.2$
ISEDBFO	$Nre = 5, Ned = 2, Nc = 10, Ns = 4, \alpha = 0.2,$ $Ped = 0.25, d_{attract} = h_{repellant} = 0.1, w_{attract} = 5, w_{repellant} = 10$
SMA	$z = 0.3, r \in [0, 1], b = [0, 1]$
MOBIFS	$Nre = 4, Ned = 2, Nc = 200, Ns = 4, \alpha = 0.2, Ped = 0.25$
BMRFO	$T(x) = x/\sqrt{1+x^2} , S = 2, r, r_1, r_2, r_3 \in [0, 1]$
IBFA	$\gamma = 1, \beta_0 = 1, \alpha = 0.5 - 0.5(t/\text{Max It})$

- (1) **Input:** missing labels and high-dimensional data and number of selected features: D
- (2) **Semisupervised learning:** labeling by Pseudocodes 2 and 3
- (3) **Initialization:** parameter setting follows Section 4.3; population initialization follows Section 3.2; iteration: $I = 0$; swimming: $s = 0$;
- (4) **While** $I < \text{max iterations}$
- (5) $I = I + 1$;
- (6) Dynamic learning by Pseudocode 4
- (7) Obtain fitness and feature subset
- (8) **If** $s < N_s // N_s$ is the number of swimming times
- (9) $s = s + 1$;
- (10) Elite population evolution by Pseudocode 5
- (11) Obtain fitness and the feature subset
- (12) **End**
- (13) **End**
- (14) **Output:** fitness and the feature subset

PSEUDOCODE 6: SBHFS.

TABLE 2: Datasets for feature selection.

Datasets	#Features	#Instances	#Class	#Smallest class	#Largest class
Australian	15	690	2	222	468
German	24	1000	2	300	700
Ionosphere	33	351	2	38	313
Colon	1999	62	2	22	40
SRBCT	2308	83	4	13	35
DLBCL	5469	77	2	25	75
Leukemia-ALLAML (LA)	7129	72	4	23	49
Central Nervous System (CNS)	7129	60	2	21	39

elimination-dispersal bacterial foraging optimization algorithm (ISEDBFO) [42], and multiobjective bacterial-inspired algorithm (MOBIFS) [48] are three recently proposed BFO variants for FS, which have good performances. ACBFO proposed an adaptive chemotaxis strategy, and ISEDBFO adopts a hyperbolic tangent function and a roulette technique to improve the search effects of BFO in FS. MOBIFS is an effective multiobjective BFO algorithm that handles FS issues using four information exchange mechanisms. The slime mold algorithm (SMA) [49], binary manta ray foraging optimization (BMRFO) [50], and improved binary butterfly algorithm (IBFA) [51] are three other bioinspired algorithms that have good performance in FS. SMA imitates slime mold's foraging behavior and introduces the composite mutation strategy and restart

strategy. BMRFO is a manta ray heuristic algorithm for FS problem solving that uses a rational transfer function. IBFA uses a new dynamic mutation operator to increase the diversity of the searching population.

Except for the abovementioned six bioinspired benchmark algorithms, to better verify the effectiveness of SBHFS, we designed two more groups of comparison experiments: comparisons with standard BFO and BCO and comparisons with semisupervised methods:

- (i) *Comparison with Standard BFO and BCO.* Based on the basic bacterial evolutionary framework, the SBHFS method has been developed with some efficient strategies. This comparison intends to evaluate the enhanced performance of the proposed

approach compared with the standard bacterial heuristic algorithm.

- (ii) *Comparisons with Semisupervised Methods.* The proposed self-adjustment, semisupervised learning method is being evaluated in this experimental group. First, original datasets are randomly divided into training and test sets. Moreover, the training set is further divided into a labeled subset and an unlabeled subset. Finally, two KNN-based semisupervised labeling techniques, semisupervised KNN (SSKNN) [38] and the best K semisupervised KNN (BKSKNN) [41], are selected to be executed on eight incompletely labeled datasets with SBHFS.

4.3. Design. In this study, all experiments were performed on a PC with Windows 10, Intel Core i7-7700, at 3.6 GHz, 8 GB RAM, and the Windows 10 operating system. Moreover, for all algorithms, the population size was set to 30, and the number of maximum iterations (\max_{iter}) was set to 100. All experiments were run independently 30 times. Due to the facility to implement KNN, this paper used KNN as the learning algorithm to assess the classification performance after FS as in literature [38, 41]. In each dataset, 70% of the samples from each class were randomly selected as the training set and the remaining 30% as the testing set. To simulate partially labeled data, this paper divided the training set into half-labeled samples and half-unlabeled samples (see Section 3.1). According to the previous experiments [6], only a small subset of tenths of the features provides the ideal solution. When the number of features for the last five datasets in Table 1 is less than 50, it is possible to attain high classification accuracy. The desired number of features (Fno.) therefore varies between 1 and 10 for the first three datasets (with reduced feature subset size) and between 5 and 50 for the remaining datasets. The parameters of all benchmark methods are given in Table 2.

For evaluation metrics (Equations.(20)–(25)), the classification error rate (denoted as Err.), true-positive rate (TPR), true-negative rate (TNR), precision (Pre), G-means (GM), and $F1$ score ($F1$) are used to assess feature selection results [52]. The effectiveness of feature selection approaches can be fully reflected by these evaluation metrics. The performance of the classification result for imbalanced data is assessed using the error rate and G-means. The TNR measures a method's capacity to isolate true-positive samples (minority samples) from all other samples, whereas the TPR measures a method's ability to isolate negative samples (majority samples) from all other samples. Precision gauges a method's capacity to distinguish genuine positive samples from all other positive samples (including true positives and false positives). A thorough evaluation of TPR and precision performance is provided by the $F1$ score.

$$Err = \frac{FP + FN}{TP + TN + FP + FN}, \quad (20)$$

$$TPR = \frac{TP}{TP + FN}, \quad (21)$$

$$TNR = \frac{TN}{TN + FP}, \quad (22)$$

$$Pre = \frac{TP}{TP + FP}, \quad (23)$$

$$GM = \sqrt{TPR + TNR}, \quad (24)$$

$$F1 = 2 * \frac{Pre * TPR}{Pre + TPR}. \quad (25)$$

Additionally, the Wilcoxon rank-sum test [53] was performed on each approach. It is marked as “=” when the p value is greater than 0.05, meaning there is no significant difference under the significance level of 5%. If the p value is less than 0.05, the recommended method is considered more significant than the comparison algorithms and marked as “+.” Otherwise, it is marked as “-.”

5. Experimental Results and Analyses

This section gives the comparison results and analyses of the three experimental groups. First, the improvement of the proposed bacterial heuristic optimization algorithm is proved by making comparisons with standard BFO and BCO for feature selection. Next, the enhanced semisupervised method is verified and discussed with two KNN-based semisupervised methods. Finally, the effectiveness of the overall proposed SBHFS method for tracking incomplete data classification is demonstrated. In Tables 3–6, the value in bold represents the best value for the current indicator. When the p value is “=,” there is no significant difference between algorithms. Therefore, the evaluation index score corresponding to the p value will not be bold.

5.1. Comparisons with Standard BFO and BCO. This comparison aims to verify the effectiveness of the proposed three strategies in BHFS, including hierarchical population initialization, dynamic learning, and elite population evolution. Table 3 shows the comparison results among the proposed bacterial heuristic optimization algorithm for FS (BHFS) and BFO for FS (BFOFS) and BCO for FS (BCOFS). The rows of Ave. and Std. show the average and standard deviation classification metrics of 30 independent runs, respectively. The rows of p show the significance values obtained by the Wilcoxon rank-sum test.

From the specific data, the feature numbers of these algorithms are consistently unchanged. This is because the controlling strategies for BHFS, BFOFS, and BCOFS are the same (see Section 4.3). Consequently, there is no difference in the significance of Fno.

On the whole, excluding Fno, BHFS obtains significantly better results in 92 out of 96 cases versus BFOFS while achieving statistically similar performance in 4 cases. Since the proposed three strategies of BHFS are the improvements of BFOFS and BCOFS, they are also the key modules that compose BHFS, where each strategy is interlinked. This result proves that BHFS is better than BFOFS and BCOFS, which reflects that our improvements are effective.

TABLE 3: The results of the comparisons with standard BFO and BCO.

Methods	Dataset	Colon							SRBCT						
		Err.	Fno.	TNR	TPR	Pre.	GM	F1	Err.	Fno.	TNR	TPR	Pre.	GM	F1
BHFS	Ave.	0.038	27.380	0.917	0.988	0.956	0.951	0.971	0.069	27.440	0.977	0.932	0.936	0.951	0.926
	Std.	0.011	15.095	0.028	0.011	0.015	0.015	0.009	0.044	15.040	0.015	0.046	0.038	0.034	0.048
BFOFS	Ave.	0.321	27.500	0.477	0.800	0.731	0.606	0.761	0.472	28.889	0.842	0.543	0.548	0.657	0.520
	Std.	0.076	15.138	0.190	0.070	0.082	0.116	0.054	0.128	15.366	0.042	0.124	0.147	0.095	0.125
BCOFS	<i>P</i>	+	=	+	+	+	+	+	+	=	+	+	+	+	+
	Ave.	0.116	27.500	0.757	0.958	0.877	0.848	0.913	0.150	27.500	0.950	0.854	0.855	0.896	0.841
	Std.	0.022	15.138	0.118	0.059	0.055	0.044	0.016	0.056	15.138	0.017	0.072	0.078	0.049	0.074
<i>P</i>	+	=	+	=	+	+	+	+	=	+	+	+	+	+	
Methods	Dataset	Ionosphere							German						
		Err.	Fno.	TNR	TPR	Pre.	GM	F1	Err.	Fno.	TNR	TPR	Pre.	GM	F1
BHFS	Ave.	0.055	12.909	0.976	0.684	0.816	0.811	0.722	0.269	10.383	0.317	0.919	0.753	0.526	0.827
	Std.	0.015	6.927	0.028	0.143	0.110	0.072	0.055	0.015	5.468	0.092	0.048	0.023	0.076	0.011
BFOFS	Ave.	0.136	15.100	0.936	0.307	0.359	0.494	0.300	0.343	11.983	0.317	0.827	0.726	0.489	0.770
	Std.	0.031	10.005	0.036	0.263	0.154	0.204	0.180	0.028	7.229	0.108	0.095	0.023	0.090	0.032
BCOFS	<i>P</i>	+	=	+	+	+	+	+	+	=	+	+	=	+	
	Ave.	0.076	16.000	0.967	0.554	0.689	0.725	0.597	0.291	12.000	0.291	0.898	0.744	0.485	0.811
	Std.	0.011	9.950	0.026	0.165	0.088	0.088	0.064	0.010	7.211	0.147	0.076	0.027	0.123	0.016
<i>P</i>	+	=	+	+	+	+	+	+	=	=	=	=	=	+	
Methods	Dataset	CNS							LA						
		Err.	Fno.	TNR	TPR	Pre.	GM	F1	Err.	Fno.	TNR	TPR	Pre.	GM	F1
BHFS	Ave.	0.093	27.420	0.837	0.942	0.925	0.884	0.931	0.000	27.440	1.000	1.000	1.000	1.000	1.000
	Std.	0.013	15.066	0.046	0.023	0.019	0.020	0.010	0.000	15.069	0.000	0.000	0.000	0.000	0.000
BFOFS	Ave.	0.428	27.500	0.350	0.723	0.690	0.482	0.703	0.136	27.500	0.886	0.853	0.945	0.867	0.895
	Std.	0.120	15.138	0.166	0.110	0.068	0.125	0.079	0.064	15.138	0.131	0.069	0.062	0.076	0.050
BCOFS	<i>P</i>	+	=	+	+	+	+	+	+	=	+	+	+	+	+
	Ave.	0.211	27.500	0.633	0.867	0.830	0.734	0.844	0.046	27.500	0.914	0.973	0.962	0.942	0.967
	Std.	0.051	15.138	0.131	0.098	0.047	0.059	0.043	0.037	15.138	0.100	0.034	0.043	0.053	0.027
<i>P</i>	+	=	+	+	+	+	+	+	=	+	+	+	+	+	
Methods	Dataset	Australian							DLBCL						
		Err.	Fno.	TNR	TPR	Pre.	GM	F1	Err.	Fno.	TNR	TPR	Pre.	GM	F1
BHFS	Ave.	0.359	5.614	0.741	0.434	0.450	0.563	0.438	0.002	27.400	0.999	0.998	1.000	0.999	0.999
	Std.	0.042	2.297	0.076	0.069	0.047	0.034	0.040	0.005	15.069	0.000	0.007	0.000	0.004	0.004
BFOFS	Ave.	0.489	7.471	0.556	0.457	0.326	0.485	0.374	0.225	27.500	0.587	0.844	0.855	0.692	0.848
	Std.	0.049	4.195	0.118	0.159	0.032	0.039	0.065	0.108	15.138	0.225	0.083	0.079	0.166	0.073
BCOFS	<i>P</i>	+	=	+	+	+	+	+	+	=	+	+	+	+	+
	Ave.	0.407	7.500	0.689	0.393	0.380	0.518	0.384	0.057	27.500	0.883	0.965	0.961	0.921	0.962
	Std.	0.037	4.183	0.065	0.052	0.038	0.029	0.033	0.041	15.138	0.112	0.041	0.037	0.062	0.028
<i>P</i>	+	=	+	+	+	+	+	+	=	+	+	+	+	+	

TABLE 4: The average computation time (minutes) of BHFS, BFOFS, and BCOFS for each run.

Datasets	Algorithms		
	BHFS	BFOFS	BCOFS
Australian	4.566	30.104	5.477
German	4.667	27.046	5.371
Ionosphere	4.019	23.034	4.182
Colon	2.420	22.398	3.360
SRBCT	2.818	25.946	3.963
DLBCL	1.803	29.433	4.388
Leukemia-ALLAML (LA)	2.833	34.795	4.282
Central Nervous System (CNS)	3.540	30.153	6.104

The bold values indicate that the SBHFS takes the least amount of time in each independent run.

From the comparison between BFOFS and BCOFS, we can see that the classification results of BCOFS perform better than those of BFOFS. This demonstrates that the improved life cycle

model in BCOFS performs better than the triple-nested loop structure, in which optimization capability is further enhanced. Moreover, compared with BFOFS and BCOFS, BHFS achieves significantly better performance in 86 out of 96 cases while obtaining statistically similar results in 10 cases. In particular, it almost achieves the best classification error rate for eight datasets. Except for the German dataset with the 26.9% classification error rate on average, BHFS for other datasets has achieved an accuracy rate of more than 90%, even the 100% accuracy rate achieved for LA and DLBCL, two microarray datasets. Thus, it is evident that BHFS outperforms both BFOFS and BCOFS. The primary reason is that BHFS has further developed the life cycle with the three proposed strategies, which improve the algorithm's search ability to locate the optimal space in the population initialization step, increase the probability of individual learning in the chemotaxis stage, and enhance the quality of population evolution in the reproduction and dispersal-elimination stages.

TABLE 5: The results of the comparisons with semisupervised methods.

Methods	Dataset	Colon							SRBCT						
		Err.	Fno.	TNR	TPR	Pre.	GM	F1	Err.	Fno.	TNR	TPR	Pre.	GM	F1
SBHFS	Ave.	0.041	27.320	0.906	0.990	0.949	0.946	0.968	0.065	27.400	0.977	0.938	0.941	0.956	0.936
	Std.	0.008	14.916	0.030	0.016	0.016	0.011	0.006	0.031	15.106	0.011	0.029	0.029	0.021	0.029
BHFS-SSKNN	Ave.	0.153	27.500	0.629	0.975	0.822	0.778	0.890	0.254	27.300	0.907	0.755	0.817	0.819	0.763
	Std.	0.046	15.138	0.138	0.040	0.054	0.082	0.030	0.046	14.818	0.018	0.051	0.044	0.042	0.052
BHFS-BKSKNN	Ave.	0.121	27.500	0.700	0.983	0.852	0.827	0.912	0.263	27.500	0.912	0.769	0.789	0.825	0.745
	Std.	0.025	15.138	0.105	0.035	0.043	0.051	0.016	0.040	15.138	0.012	0.040	0.068	0.030	0.052
	<i>P</i>	+	=	+	=	+	+	+	+	=	+	+	+	+	+
Methods	Dataset	Ionosphere							German						
		Err.	Fno.	TNR	TPR	Pre.	GM	F1	Err.	Fno.	TNR	TPR	Pre.	GM	F1
SBHFS	Ave.	0.050	13.655	0.994	0.582	0.926	0.756	0.704	0.273	10.517	0.294	0.913	0.753	0.506	0.824
	Std.	0.013	7.794	0.005	0.110	0.049	0.072	0.090	0.016	5.735	0.112	0.028	0.025	0.090	0.005
BHFS-SSKNN	Ave.	0.152	15.500	0.860	0.745	0.384	0.798	0.505	0.282	11.250	0.236	0.933	0.737	0.458	0.822
	Std.	0.016	7.634	0.020	0.112	0.035	0.055	0.048	0.010	6.398	0.095	0.049	0.019	0.076	0.011
BHFS-BKSKNN	Ave.	0.063	14.182	0.982	0.554	0.783	0.734	0.644	0.289	11.583	0.219	0.921	0.734	0.441	0.817
	Std.	0.012	8.280	0.008	0.103	0.082	0.069	0.082	0.008	7.128	0.089	0.036	0.016	0.076	0.007
	<i>P</i>	+	=	+	=	+	=	=	+	=	=	=	+	=	+
Methods	Dataset	CNS							LA						
		Err.	Fno.	TNR	TPR	Pre.	GM	F1	Err.	Fno.	TNR	TPR	Pre.	GM	F1
SBHFS	Ave.	0.104	27.380	0.740	0.973	0.885	0.846	0.925	0.012	27.500	0.989	0.988	0.995	0.988	0.991
	Std.	0.006	15.033	0.044	0.024	0.017	0.015	0.005	0.011	15.138	0.020	0.010	0.009	0.013	0.008
BHFS-SSKNN	Ave.	0.133	27.500	0.617	0.992	0.839	0.780	0.909	0.114	27.500	0.671	0.987	0.869	0.809	0.923
	Std.	0.029	15.138	0.081	0.026	0.027	0.052	0.019	0.044	15.138	0.151	0.028	0.058	0.082	0.029
BHFS-BKSKNN	Ave.	0.133	27.400	0.617	0.992	0.840	0.779	0.909	0.032	27.500	1.000	0.953	1.000	0.976	0.976
	Std.	0.039	14.976	0.112	0.026	0.040	0.070	0.025	0.022	15.138	0.000	0.032	0.000	0.016	0.017
	<i>P</i>	+	=	+	-	+	+	=	=	=	=	=	=	=	=
Methods	Dataset	Australian							DLBCL						
		Err.	Fno.	TNR	TPR	Pre.	GM	F1	Err.	Fno.	TNR	TPR	Pre.	GM	F1
SBHFS	Ave.	0.350	5.871	0.771	0.398	0.454	0.553	0.423	0.003	27.480	0.997	0.996	0.999	0.996	0.998
	Std.	0.016	2.600	0.017	0.043	0.029	0.030	0.036	0.006	15.105	0.011	0.008	0.004	0.006	0.004
BHFS-SSKNN	Ave.	0.483	7.600	0.419	0.724	0.351	0.489	0.462	0.117	27.500	0.550	1.000	0.863	0.740	0.927
	Std.	0.058	2.716	0.205	0.259	0.081	0.135	0.154	0.021	15.138	0.081	0.000	0.022	0.053	0.012
BHFS-BKSKNN	Ave.	0.330	6.308	0.788	0.421	0.488	0.575	0.451	0.052	27.500	0.800	1.000	0.936	0.892	0.966
	Std.	0.020	2.594	0.021	0.050	0.039	0.035	0.042	0.034	15.138	0.131	0.000	0.041	0.073	0.022
	<i>P</i>	-	=	=	=	+	+	+	+	=	+	=	+	+	+

The bold values for each method mean that they achieve the best results under the evaluation index.

Table 4 illustrates the average calculation time for feature selection and classification in each run. Compared to all bacterial-based methods, BHFS achieves a superior classification effect with less computational complexity. Throughout the iteration period, the computing time of the BFO algorithm increases exponentially due to its nested structure. However, life cycle enhancement offered by BCO streamlines this procedure, hence reducing computing cost dramatically. Inspired by BCO, BHFS modifies the parts of the population update based on BCO so that reproduction and dispersal-elimination operations can be carried out just one at a time, and the algorithm is additionally programmed with a rule to instantly stop iterating when the ideal solution occurs repeatedly.

5.2. Comparisons with Semisupervised Methods. Since BHFS demonstrates its superiority and usefulness in Section 5.1, this section evaluates the effectiveness of the proposed self-

adjustment, semisupervised KNN strategy for BHFS. In the following context, we refer to BHFS with the self-adjusted, semisupervised KNN strategy as SBHFS. The compared two semisupervised techniques based on KNN are as follows: One is the semisupervised KNN (SSKNN) [38], which assigns the unlabeled sample's label to the label of the labeled sample that is closest to it. The best K semisupervised KNN (BKSKNN) [41] is another comparative technique. By learning about their neighbors, BKSKNN also labels unlabeled samples. In contrast, BKSKNN has two steps as opposed to SSKNN, the first of which is to compute the accuracy of the labeling result of KNN using various K values. K is then set from 1 to 51. The process finds the best K value with the highest level of labeling accuracy and then uses the best-labeled data to perform the subsequent procedure. These two semisupervised learning approaches are embedded into BHFS for the comparison of the effectiveness

TABLE 6: The result of the comparisons with benchmark methods.

Methods	Indexes	Colon							SRBCT						
		Err.	Fno.	TNR	TPR	Pre.	GM	F1	Err.	Fno.	TNR	TPR	Pre.	GM	F1
SBHFS	Ave.	0.041	27.5	0.92	0.982	0.957	0.949	0.968	0.064	27.5	0.979	0.937	0.94	0.956	0.932
	Std.	0.015	15.138	0.038	0.015	0.021	0.02	0.012	0.036	15.138	0.012	0.042	0.033	0.029	0.044
	P	+	+	+	+	+	+	+	+	+	+	+	+	+	+
ACBFO	Ave.	0.176	390.6	0.665	0.91	0.84	0.769	0.871	0.159	458.35	0.945	0.863	0.877	0.894	0.853
	Std.	0.028	15.364	0.055	0.031	0.026	0.035	0.022	0.03	19.594	0.011	0.028	0.025	0.025	0.028
	P	+	+	+	+	+	+	+	+	+	+	+	+	+	+
ISEDBFO	Ave.	0.153	232.45	0.679	0.938	0.848	0.791	0.888	0.162	246.4	0.944	0.859	0.869	0.891	0.842
	Std.	0.029	56.974	0.073	0.029	0.03	0.046	0.021	0.043	73.685	0.015	0.043	0.04	0.034	0.048
	P	+	+	+	+	+	+	+	+	+	+	+	+	+	+
SMA	Ave.	0.135	4.3	0.773	0.916	0.887	0.832	0.896	0.134	7.833	0.955	0.872	0.887	0.907	0.863
	Std.	0.035	1.302	0.074	0.06	0.031	0.039	0.031	0.035	5.834	0.011	0.034	0.038	0.025	0.038
	P	+	-	+	+	+	+	+	+	-	+	+	+	+	+
MOBIFS	Ave.	0.25	19.166	0.438	0.79	0.91	0.478	0.845	0.48	18.924	0.914	0.471	0.731	0.561	0.62
	Std.	0.108	2.104	0.489	0.065	0.069	0.402	0.067	0.117	3.621	0.083	0.073	0.326	0.218	0.079
	P	+	=	=	+	=	+	+	+	=	=	+	+	+	+
BMRFO	Ave.	0.179	34.75	0.666	0.905	0.838	0.765	0.866	0.129	216	0.955	0.888	0.897	0.914	0.874
	Std.	0.032	35.84	0.092	0.033	0.038	0.055	0.023	0.048	146.98	0.017	0.044	0.035	0.035	0.049
	P	+	=	+	+	+	+	+	+	+	-	+	+	+	+
IBFA	Ave.	0.221	988.3	0.571	0.886	0.801	0.695	0.837	0.19	1140.8	0.93	0.84	0.86	0.88	0.83
	Std.	0.022	20.63	0.04	0.027	0.018	0.029	0.019	0.02	21.71	0.01	0.02	0.02	0.02	0.02
	P	+	+	+	+	+	+	+	+	+	+	+	+	+	+

Methods	Indexes	Ionosphere							German						
		Err.	Fno.	TNR	TPR	Pre.	GM	F1	Err.	Fno.	TNR	TPR	Pre.	GM	F1
SBHFS	Ave.	0.062	14.491	0.963	0.721	0.750	0.829	0.716	0.246	10.967	0.407	0.917	0.777	0.595	0.839
	Std.	0.03	8.114	0.042	0.099	0.128	0.043	0.062	0.024	5.954	0.125	0.039	0.036	0.114	0.008
	P	+	-	=	=	-	-	-	-	+	-	+	+	+	+
ACBFO	Ave.	0.104	4.85	0.976	0.754	0.949	0.855	0.834	0.288	4.5	0.351	0.867	0.758	0.542	0.807
	Std.	0.011	0.988	0.013	0.038	0.027	0.019	0.021	0.01	1.539	0.064	0.027	0.013	0.049	0.009
	P	+	-	=	=	-	-	-	-	+	-	+	+	+	+
ISEDBFO	Ave.	0.085	3.05	0.97	0.817	0.94	0.89	0.873	0.275	4.1	0.409	0.86	0.774	0.588	0.814
	Std.	0.008	0.51	0.009	0.022	0.017	0.012	0.014	0.01	0.553	0.059	0.021	0.014	0.041	0.007
	P	+	-	=	-	-	-	-	+	-	=	+	=	+	+
SMA	Ave.	0.091	2.5	0.961	0.816	0.925	0.884	0.864	0.305	2.9	0.332	0.851	0.749	0.513	0.793
	Std.	0.013	0.607	0.015	0.021	0.025	0.014	0.019	0.04	0.852	0.083	0.039	0.028	0.082	0.032
	P	+	-	+	-	-	-	-	+	-	+	+	+	+	+
MOBIFS	Ave.	0.677	11.732	0.881	0.118	0.738	0.321	0.203	0.29	9.252	0.497	0.737	0.924	0.601	0.82
	Std.	0.059	2.939	0.052	0.016	0.064	0.031	0.026	0.027	3.434	0.125	0.009	0.038	0.08	0.02
	P	+	=	+	+	=	+	+	+	=	=	+	-	=	+
BMRFO	Ave.	0.09	3.1	0.965	0.811	0.936	0.883	0.865	0.279	4.6	0.342	0.884	0.76	0.537	0.816
	Std.	0.009	0.968	0.012	0.023	0.02	0.012	0.014	0.015	1.875	0.101	0.038	0.021	0.087	0.011
	P	+	-	=	-	-	-	-	+	-	+	+	+	+	+
IBFA	Ave.	0.116	9.7	0.978	0.716	0.951	0.835	0.814	0.258	9.25	0.441	0.871	0.785	0.617	0.825
	Std.	0.009	1.809	0.009	0.02	0.019	0.013	0.016	0.016	1.743	0.042	0.011	0.013	0.031	0.01
	P	+	=	=	=	-	=	-	+	=	=	+	=	=	+

Methods	Indexes	CNS							LA						
		Err.	Fno.	TNR	TPR	Pre.	GM	F1	Err.	Fno.	TNR	TPR	Pre.	GM	F1
SBHFS	Ave.	0.082	27.5	0.82	0.967	0.918	0.887	0.94	0.002	27.48	1	0.997	1	0.999	0.999
	Std.	0.018	15.138	0.042	0.022	0.018	0.023	0.014	0.004	15.127	0	0.006	0	0.003	0.003
	P	+	+	+	+	+	+	+	+	+	+	+	+	+	+
ACBFO	Ave.	0.38	1411.2	0.444	0.716	0.711	0.53	0.702	0.119	1394.7	0.942	0.853	0.970	0.894	0.905
	Std.	0.028	37.047	0.085	0.04	0.035	0.061	0.025	0.011	35.534	0.029	0.011	0.012	0.016	0.009
	P	+	+	+	+	+	+	+	+	+	+	+	+	+	+
ISEDBFO	Ave.	0.38	731.8	0.487	0.695	0.724	0.495	0.689	0.121	647.3	0.936	0.858	0.969	0.892	0.905
	Std.	0.031	108.87	0.071	0.059	0.021	0.059	0.038	0.019	30.78	0.046	0.006	0.021	0.027	0.013
	P	+	+	+	+	+	+	+	+	+	+	+	+	+	+
SMA	Ave.	0.031	108.87	0.071	0.059	0.021	0.059	0.038	0.036	3.6	0.969	0.963	0.986	0.964	0.973
	Std.	0.067	1.353	0.112	0.07	0.052	0.093	0.056	0.028	2.683	0.047	0.034	0.021	0.031	0.021
	P	-	+	+	+	+	+	+	+	-	+	+	+	+	+
MOBIFS	Ave.	0.333	19.262	0.722	0.683	0.901	0.682	0.775	0.174	20.734	0.902	0.815	0.981	0.854	0.889
	Std.	0.144	2.305	0.357	0.101	0.089	0.285	0.095	0.094	4.515	0.18	0.082	0.029	0.123	0.057
	P	+	=	=	=	=	=	=	+	=	+	+	+	+	+

TABLE 6: Continued.

Methods	Indexes	Colon							SRBCT						
		Err.	Fno.	TNR	TPR	Pre.	GM	F1	Err.	Fno.	TNR	TPR	Pre.	GM	F1
BMRFO	Ave.	0.38	358.3	0.458	0.701	0.717	0.518	0.699	0.081	16.35	0.925	0.916	0.967	0.917	0.938
	Std.	0.059	470.91	0.103	0.061	0.056	0.096	0.051	0.029	17.236	0.063	0.033	0.027	0.038	0.022
	<i>P</i>	+	+	+	+	+	+	+	+	+	+	+	+	+	+
IBFA	Ave.	0.404	3517.7	0.413	0.694	0.689	0.518	0.69	0.114	3492.3	0.948	0.86	0.972	0.901	0.91
	Std.	0.025	36.131	0.049	0.029	0.021	0.042	0.021	0.012	23.086	0.019	0.013	0.009	0.013	0.009
	<i>P</i>	+	+	+	+	+	+	+	+	+	+	+	+	+	+
Methods	Indexes	Australian							DLBCL						
		Err.	Fno.	TNR	TPR	Pre.	GM	F1	Err.	Fno.	TNR	TPR	Pre.	GM	F1
SBHFS	Ave.	0.357	5.857	0.7	0.524	0.456	0.605	0.487	0	27.46	1	1	1	1	1
	Std.	0.024	2.474	0.03	0.034	0.032	0.025	0.03	0	15.124	0	0	0	0	0
	<i>P</i>	+	+	+	+	+	+	+	+	+	+	+	+	+	+
ACBFO	Ave.	0.145	1	0.925	0.799	0.931	0.859	0.859	0.108	1084.6	0.798	0.924	0.934	0.852	0.927
	Std.	0	0	0	0	0	0	0	0.018	30.659	0.063	0.022	0.017	0.033	0.012
	<i>P</i>	-	-	-	-	-	-	-	+	+	+	+	+	+	+
ISEDBFO	Ave.	0.145	1	0.925	0.799	0.931	0.859	0.859	0.057	623.6	0.887	0.962	0.964	0.92	0.962
	Std.	0	0	0	0	0	0	0	0.015	127.109	0.048	0.017	0.015	0.025	0.01
	<i>P</i>	-	-	-	-	-	-	-	+	+	+	+	+	+	+
SMA	Ave.	0.145	1	0.925	0.799	0.93	0.859	0.859	0.097	12.55	0.81	0.936	0.939	0.864	0.935
	Std.	0	0	0	0	0	0	0	0.036	17.111	0.096	0.028	0.029	0.06	0.024
	<i>P</i>	-	-	-	-	-	-	-	+	-	+	+	+	+	+
MOBIFS	Ave.	0.448	2.556	0.734	0.393	0.625	0.537	0.481	0.22	19.633	0.937	0.784	0.966	0.874	0.865
	Std.	0.024	1.262	0.002	0.015	0.029	0.01	0.005	0.097	2.387	0.081	0.071	0.041	0.068	0.058
	<i>P</i>	-	=	-	-	-	-	=	+	=	+	+	+	+	+
BMRFO	Ave.	0.144	1.5	0.912	0.811	0.922	0.859	0.861	0.104	76.05	0.75	0.942	0.927	0.817	0.931
	Std.	0.002	1.235	0.031	0.028	0.021	0.001	0.006	0.039	81.127	0.112	0.042	0.031	0.089	0.028
	<i>P</i>	-	-	-	-	-	-	-	+	+	+	+	+	+	+
IBFA	Ave.	0.154	2.2	0.87	0.826	0.893	0.845	0.855	0.09	2684.9	0.81	0.94	0.94	0.87	0.94
	Std.	0.014	1.152	0.062	0.027	0.04	0.02	0.008	0.02	22.44	0.04	0.01	0.01	0.02	0.01
	<i>P</i>	-	=	=	-	=	=	-	+	+	+	+	+	+	+

The values that are in bold show that each method got the best results based on the evaluation metrics.

of different semisupervised methods, and they are recorded as BHFS-SSKNN and BHFS-BKSKNN, respectively. Table 5 shows the average, standard deviation classification metric, and statistical test results of different semisupervised learning approaches for benchmark datasets. Since all three methods are based on BHFS, their feature subset size control methodologies are identical (see Section 4.3). Consequently, the significance of Fno. does not change in the three methods.

SBHFS outperforms BHFS-SSKNN and BHFS-BKSKNN in the majority of classification evaluation metrics, demonstrating the efficacy of the self-adjusted, semisupervised KNN technique. Self-adjusted, semisupervised KNN will adaptively update the K value to find a better label for each sample from varying data sizes, whereas SSKNN will simply apply the fixed K value that limits the algorithm's performance.

Compared with BHFS-SSKNN, SBHFS obtains a lower error rate in all data cases. In other classification metrics, BHFS-SSKNN shows better performance with the true-positive rate (TPR), indicating that the SSKNN method is more capable of correctly labeling positive samples, while SBHFS can achieve significantly better or similar performance with the true-negative rate (TNR). This proves that

using TPR or TNR metrics alone to judge the performance of algorithms is one-sided. Therefore, it is necessary to deeply analyze the algorithm effect through the remaining three comprehensive evaluation indicators (Pre, GM, and $F1$). The results demonstrate that the Pre scores of SBHFS are 100 percent superior to those of BHFS-SSKNN, while the GM scores of SBHFS are higher for five out of eight datasets, and the $F1$ scores are higher for half of the datasets. Thus, self-adjusted, semisupervised KNN does improve the performance of SBHFS.

Figure 4 shows the bar chart of the comparison results with semisupervised methods. The horizontal axis corresponds to evaluation metrics. Fno. is excluded since the comparison algorithms use the same feature subset size control methods. The ordinate represents each algorithm's score, with larger values indicating superior performance. From Table 5 and Figure 4, we can see that SBHFS can achieve statistically significant better classification performance for all high-dimensional datasets with the most different metrics. For benchmark datasets, statistical significance is not as obvious. One reason is that in a low-dimensional space, the KNN-based semisupervised learning method is less affected by the value of K . Therefore, we conclude that the proposed self-adjusted, semisupervised

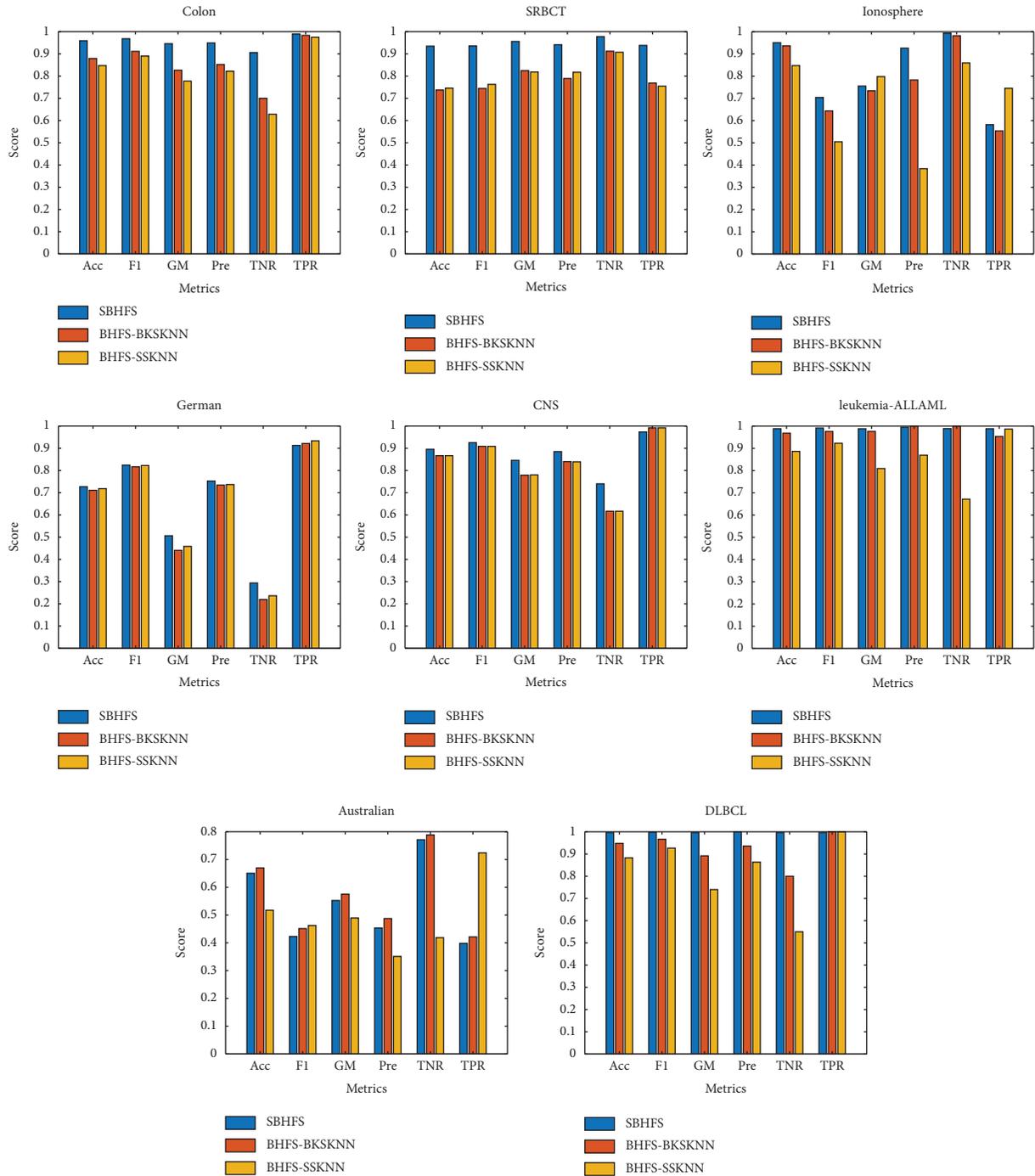


FIGURE 4: The bar chart of the comparison results with semisupervised methods.

KNN mechanism is more suitable for high-dimensional data analysis.

5.3. *Comparisons with Benchmark Methods.* We further compare the proposed SBHFS method with other bio-inspired wrapper FS algorithms. Table 6 shows different evaluation metrics (i.e., Err., Fno., TPR, TNR, Pre, GM, and F1) of SBHFS and benchmark methods for test sets. In general, SBHFS achieves competitive results compared to the

other six bioinspired feature selection methods, which means that SBHFS is superior to other bioinspired wrapper FS algorithms.

From the specific results, SBHFS performs best for five datasets (i.e., LA, CNS, Colon, SRBCT, and DLBCL) with the most evaluation metrics. However, for the Australian dataset, the proposed method does not perform best. Comparing the results of other algorithms reveals that the effect of SBHFS may be influenced by the sparsity of features. To be specific, according to the results for the Australian

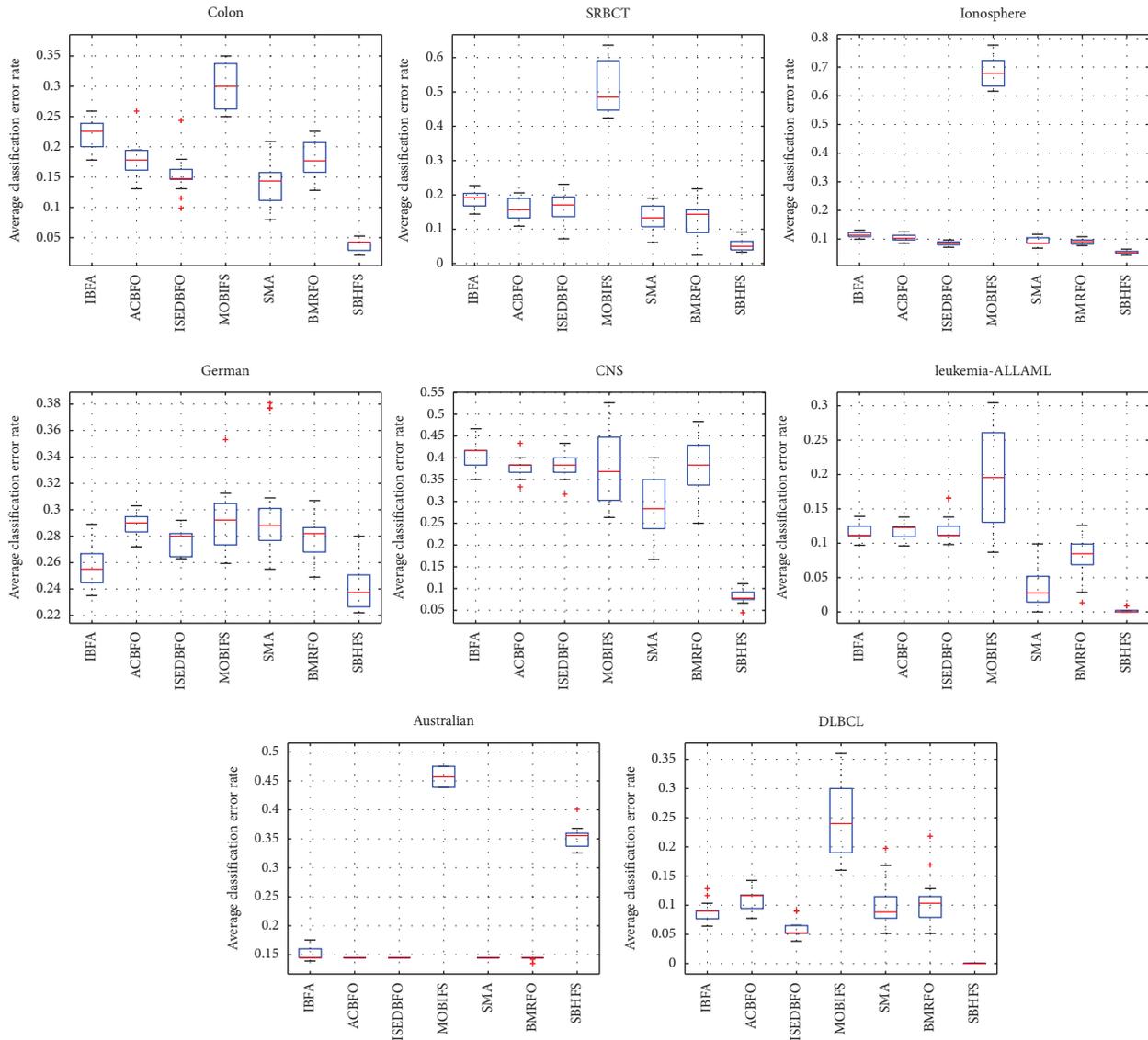


FIGURE 5: The box-line plot of comparison results with benchmark methods.

dataset, we can see that smaller subsets produce better results and that the number of effective features of the Australian dataset is about 1 or 2. Nevertheless, SBHFS sets the subset size to integers between 1 and 10 when the total dataset contains fewer than 50 features (see Section 4.3). This increases the average size of the subset in each iteration, which exceeds the effective feature size of the Australian dataset. Except for this, based on the results of statistical significance tests for all datasets, SBHFS achieves considerably enhanced efficiency in 229 of 336 cases (39 cases with the “=” p value are excluded), which illustrates that SBHFS performs well for most of the datasets, especially for high-dimensional ones.

From the perspective of the fundamental algorithm, SBHFS achieves notable significance in 150 out of 224 cases (25 cases with the “=” p value are excluded) in comparison to four other bacterial-based FS methods, i.e., ACBFO, ISEDBFO, SMA, and MOBIFS. The results demonstrate that

the improvements in SBHFS are better than in other bacterial-based algorithms in this study. The proposed strategies optimize the searching ability of the algorithm, achieving smaller classification error rates and better results on other evaluation indexes. Moreover, the superiority of SBHFS is more obvious for high-dimensional datasets (i.e., #features > 1999). For example, compared with ISEDBFO, SBHFS achieves 11.2% lower Err. for Colon (#features=1999) and 9.8% lower Err. for SRBCT (#features=2308). This demonstrates that the dimension redundancy capability of the suggested feature selection approach is satisfactory.

Furthermore, compared with BMRFO and IBFA, SBHFS achieves superior performance in 79 out of 112 cases (14 cases with the “=” p value are excluded). Specifically, compared to IBFA, SBHFS achieves 9.9% higher $F1$ for LA (#Features=7129) and 1.4% higher $F1$ for SRBCT (#

Features = 2308). Compared to bacterial-based algorithms, the effects between SBHFS and other bioinspired algorithms are better. This shows that the proposed modified bacterial-based FS algorithm is better for solving dimension reduction problems and that SBHFS can be used not only for high-dimensional datasets but also for some low-dimensional datasets.

To verify the stability of SBHFS, the boxplot in Figure 5 shows the comparison results of the SBHFS with other bioinspired FS methods. According to the boxplot, except for the Australian dataset with the 35.7% average classification error rate, SBHFS has achieved the best accuracy rate compared with other algorithms for other datasets. Moreover, the median results show that SBHFS generally achieves lower error rates, and the width of boxes indicates that SBHFS is more stable than other comparison methods. This is due to the dynamic learning method that allows the bacterial population of each iteration to move closer to the optimal solution instead of searching randomly.

6. Conclusions

This paper presents a semisupervised bacterial heuristic feature selection algorithm (SBHFS) to address label incomplete and high-dimensional classification problems. The self-adjusted, semisupervised KNN strategy can reconstruct labels effectively with the help of the two-step self-training mechanism, and the improved bacterial heuristic method can enhance the searching precision by increasing feature selection variety and cooperating with hierarchical population initialization, dynamic learning, and elite population evolution strategies. To be specific, hierarchical population initialization accelerates the convergence of the algorithm with the help of the SU feature ranking method and the proposed layer mechanism. Then, the dynamic learning strategy increases the diversity of the feature subset because it promotes the communication of searching bacteria. Furthermore, the proposed elite population evolution strategy changes the population update method of the bacterial-based algorithm and improves its optimization performance. The comparisons with the semisupervised methods show that the proposed semisupervised learning method is effective for labeling incomplete data, especially for high-dimensional datasets.

Although the proposed SBHFS approach has shown promise in high-dimensional classification with missing labels, the proposed semisupervised approach is based on the enhancement of the KNN semisupervised technique, and the semisupervised method based on other learners is not considered. This may limit the efficiency of the bacterial heuristic algorithm in FS for classification issues involving plenty of sparse features. Considering this information in feature selection may help bacterial heuristic algorithms achieve better results, although this is challenging to accomplish. In our future endeavors, we will consider this direction.

Data Availability

All relevant data are included within the paper.

Ethical Approval

Ethical approval was not required for this study.

Conflicts of Interest

The authors declare that there are no conflicts of interest regarding the publication of this paper.

Authors' Contributions

Hong Wang conceptualized the study. Hong Wang, Yikun Ou, and Lijing Tan designed the methodology. Hong Wang, Yikun Ou, Yixin Wang, Tongtong Xing, and Lijing Tan validated the data. Hong Wang, Yixin Wang, and Tongtong Xing carried out the formal analysis. Hong Wang, Yikun Ou, Yixin Wang, Tongtong Xing, and Lijing Tan wrote and prepared the original draft. Hong Wang and Lijing Tan wrote and reviewed and edited the manuscript. Hong Wang and Lijing Tan were involved in funding acquisition. All authors read and agreed to the published version of the manuscript.

Acknowledgments

This work was supported by the National Natural Science Foundation of China (Grant no. 71901152), the Natural Science Foundation of Guangdong Province (Grant no. 2020A1515010752), and the Guangdong Innovation Team Project of Intelligent Management and Cross Innovation (2021WCXTD002).

References

- [1] W. Wang, X. Yang, X. Li, and J. Tang, "Convolutional-capsule network for gastrointestinal endoscopy image classification," *International Journal of Intelligent Systems*, vol. 37, no. 9, pp. 5796–5815, 2022.
- [2] D. Zhang, L. Yao, K. Chen, S. Wang, X. Chang, and Y. Liu, "Making sense of spatio-temporal preserving representations for eeg-based human intention recognition," *IEEE Transactions on Cybernetics*, vol. 50, no. 7, pp. 3033–3044, 2020.
- [3] P. Wang, B. Xue, J. Liang, and M. Zhang, "Multiobjective differential evolution for feature selection in classification," *IEEE Transactions on Cybernetics*, pp. 1–15, 2021.
- [4] J. Piri and P. Mohapatra, "An analytical study of modified multi-objective Harris hawk optimizer towards medical data feature selection," *Computers in Biology and Medicine*, vol. 135, Article ID 104558, 2021.
- [5] S. Sreejith, H. Khanna Nehemiah, and A. Kannan, "Clinical data classification using an enhanced smote and chaotic evolutionary feature selection," *Computers in Biology and Medicine*, vol. 126, Article ID 103991, 2020.
- [6] H. Wang, X. Jing, and B. Niu, "A discrete bacterial algorithm for feature selection in classification of microarray gene expression cancer data," *Knowledge-Based Systems*, vol. 126, pp. 8–19, 2017.
- [7] M. Liu, J. Zhang, C. Lian, and D. Shen, "Weakly supervised deep learning for brain disease prognosis using mri and incomplete clinical scores," *IEEE Transactions on Cybernetics*, vol. 50, no. 7, pp. 3381–3392, 2020.

- [8] Z. Chen, Y. Liu, Y. Zhang, R. Jin, J. Tao, and L. Chen, "Low-rank sparse feature selection with incomplete labels for alzheimer's disease progression prediction," *Computers in Biology and Medicine*, vol. 147, Article ID 105705, 2022.
- [9] H. Liu and H. Motoda, *Feature Extraction, Construction and Selection: A Data Mining Perspective*, Springer Science & Business Media, Assinippi Park Norwell, MA USA, 1998.
- [10] M. M. Ahmed and T. Palaniswamy, "A novel TMGWO-SLBNC-based multidimensional feature subset selection and classification framework for frequent diagnosis of breast lesion abnormalities," *International Journal of Intelligent Systems*, vol. 37, no. 3, pp. 2131–2162, 2021.
- [11] S. Liu, H. Wang, W. Peng, and W. Yao, "A surrogate-assisted evolutionary feature selection algorithm with parallel random grouping for high-dimensional classification," *IEEE Transactions on Evolutionary Computation*, 2022.
- [12] S. Jadhav, H. He, and K. Jenkins, "Information gain directed genetic algorithm wrapper feature selection for credit rating," *Applied Soft Computing*, vol. 69, pp. 541–553, 2018.
- [13] L. Sun, T. Wang, W. Ding, J. Xu, and Y. Lin, "Feature selection using Fisher score and multilabel neighborhood rough sets for multilabel classification," *Information Sciences*, vol. 578, pp. 887–912, 2021.
- [14] N. Maleki, Y. Zeinali, and S. T. A. Niaki, "A k-NN method for lung cancer prognosis with the use of a genetic algorithm for feature selection," *Expert Systems with Applications*, vol. 164, Article ID 113981, 2021.
- [15] O. Gokalp, E. Tasci, and A. Ugur, "A novel wrapper feature selection algorithm based on iterated greedy metaheuristic for sentiment classification," *Expert Systems with Applications*, vol. 146, Article ID 113176, 2020.
- [16] F. Dornaika and A. Khoder, "Linear embedding by joint robust discriminant analysis and inter-class sparsity," *Neural Networks*, vol. 127, pp. 141–159, 2020.
- [17] E. Hancer, B. Xue, and M. Zhang, "A survey on feature selection approaches for clustering," *Artificial Intelligence Review*, vol. 53, no. 6, pp. 4519–4545, 2020.
- [18] J. Fan and R. Li, "Comment: feature screening and variable selection via iterative ridge regression," *Technometrics*, vol. 62, no. 4, pp. 434–437, 2020.
- [19] X. Song, M. T. Liu, Q. Liu, and B. Niu, "Hydrological cycling optimization-based multiobjective feature-selection method for customer segmentation," *International Journal of Intelligent Systems*, vol. 36, no. 5, pp. 2347–2366, 2021.
- [20] H. Liu, M. Zhou, and Q. Liu, "An embedded feature selection method for imbalanced data classification," *IEEE/CAA Journal of Automatica Sinica*, vol. 6, no. 3, pp. 703–715, 2019.
- [21] B. Tran, B. Xue, and M. Zhang, "Variable-length particle swarm optimization for feature selection on high-dimensional classification," *IEEE Transactions on Evolutionary Computation*, vol. 23, no. 3, pp. 473–487, 2019.
- [22] X.-F. Song, Y. Zhang, Y.-N. Guo, X.-Y. Sun, and Y.-L. Wang, "Variable-size cooperative coevolutionary particle swarm optimization for feature selection on high-dimensional data," *IEEE Transactions on Evolutionary Computation*, vol. 24, no. 5, pp. 882–895, 2020.
- [23] J. González, J. Ortega, M. Damas, P. Martín Smith, and J. Q. Gan, "A new multi-objective wrapper method for feature selection – accuracy and stability analysis for bci," *Neurocomputing*, vol. 333, pp. 407–418, 2019.
- [24] T. Gangavarapu and N. Patil, "A novel filter-wrapper hybrid greedy ensemble approach optimized using the genetic algorithm to reduce the dimensionality of high-dimensional biomedical datasets," *Applied Soft Computing*, p. 81, 2019.
- [25] B. Xue, M. Zhang, and W. N. Browne, "Particle swarm optimization for feature selection in classification: a multi-objective approach," *IEEE Transactions on Cybernetics*, vol. 43, no. 6, pp. 1656–1671, 2013.
- [26] R. N. Khushaba, A. Al Ani, and A. Al Jumaily, "Feature subset selection using differential evolution and a statistical repair mechanism," *Expert Systems with Applications*, vol. 38, no. 9, pp. 11515–11526, 2011.
- [27] C. F. Tsai, W. Eberle, and C. Y. Chu, "Genetic algorithms in feature and instance selection," *Knowledge-Based Systems*, vol. 39, pp. 240–247, 2013.
- [28] H. Wang, L. Tan, and B. Niu, "Feature selection for classification of microarray gene expression cancers using bacterial colony optimization with multi-dimensional population," *Swarm and Evolutionary Computation*, vol. 48, pp. 172–181, 2019.
- [29] K. M. Passino, "Biomimicry of bacterial foraging for distributed optimization and control," *IEEE Control Systems Magazine*, vol. 22, no. 3, pp. 52–67, 2002.
- [30] H. Wang, B. Niu, and L. Tan, "Bacterial colony algorithm with adaptive attribute learning strategy for feature selection in classification of customers for personalized recommendation," *Neurocomputing*, vol. 452, pp. 747–755, 2021.
- [31] B. Niu and H. Wang, "Bacterial colony optimization," *Discrete Dynamics in Nature and Society*, vol. 2012, Article ID 698057, 28 pages, 2012.
- [32] H. Wang and B. Niu, "A novel bacterial algorithm with randomness control for feature selection in classification," *Neurocomputing*, vol. 228, pp. 176–186, 2017.
- [33] Q. Q. Pang and L. Zhang, "Semi-supervised Neighborhood Discrimination index for Feature Selection," *Knowledge-Based Systems*, vol. 204, 2020.
- [34] W. Yang, K. Xia, T. Li, M. Xie, and F. Song, "A multi-strategy marine predator algorithm and its application in joint regularization semi-supervised elm," *Mathematics*, vol. 9, no. 3, p. 291, 2021.
- [35] Y. Zhang, H. G. Li, Q. Wang, and C. Peng, "A filter-based bare-bone particle swarm optimization algorithm for unsupervised feature selection," *Applied Intelligence*, vol. 49, no. 8, pp. 2889–2898, 2019.
- [36] K. Chen, L. Yao, D. Zhang, X. Wang, X. Chang, and F. Nie, "A semisupervised recurrent convolutional attention model for human activity recognition," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 5, pp. 1747–1756, 2020.
- [37] M. Luo, X. Chang, L. Nie, Y. Yang, A. G. Hauptmann, and Q. Zheng, "An adaptive semisupervised feature analysis for video semantic recognition," *IEEE Transactions on Cybernetics*, vol. 48, no. 2, pp. 648–660, 2018.
- [38] D. Zhang, L. Jiao, X. Bai, S. Wang, and B. Hou, "A robust semi-supervised svm via ensemble learning," *Applied Soft Computing*, vol. 65, pp. 632–643, 2018.
- [39] K. Liu, X. Yang, H. Yu, J. Mi, P. Wang, and X. Chen, "Rough set based semi-supervised feature selection via ensemble selector," *Knowledge-Based Systems*, vol. 165, pp. 282–296, 2019.
- [40] A. Calma, T. Reitmaier, and B. Sick, "Resp-knn: A probabilistic K-nearest neighbor classifier for sparsely labeled data," in *Proceedings of the International Joint Conference On Neural Networks*, pp. 4040–4047, BC, Canada, November 2016.
- [41] K. Mehta, A. Jain, J. Mangalagiri, S. Menon, P. Nguyen, and D. R. Chapman, "Lung nodule classification using biomarkers, volumetric radiomics, and 3d cnns," *Journal of Digital Imaging*, vol. 34, no. 3, pp. 647–666, 2021.

- [42] Y. P. Chen, Y. Li, G. Wang et al., "A novel bacterial foraging optimization algorithm for feature selection," *Expert Systems with Applications*, vol. 83, pp. 1–17, 2017.
- [43] M. Wang and H. Chen, "Chaotic multi-swarm whale optimizer boosted support vector machine for medical diagnosis," *Applied Soft Computing*, p. 88, 2020.
- [44] Z. Zeng, L. Guan, W. Zhu, J. Dong, and J. Li, "Face recognition based on svm optimized by the improved bacterial foraging optimization algorithm," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 33, no. 7, Article ID 1956007, 2019.
- [45] M. Kaur and S. Kadam, "A novel multi-objective bacteria foraging optimization algorithm (MOBFOA) for multi-objective scheduling," *Applied Soft Computing*, vol. 66, pp. 183–195, 2018.
- [46] J. Too, A. R. Abdullah, and N. Mohd Saad, "A new co-evolution binary particle swarm optimization with multiple inertia weight strategy for feature selection," *Informatics*, vol. 6, no. 2, p. 21, 2019.
- [47] W. Pei, B. Xue, L. Shang, and M. Zhang, "Genetic programming for development of cost-sensitive classifiers for binary high-dimensional unbalanced classification," *Applied Soft Computing*, vol. 101, Article ID 106989, 2021.
- [48] B. Niu, W. Yi, L. Tan, S. Geng, and H. Wang, "A multi-objective feature selection method based on bacterial foraging optimization," *Natural Computing*, vol. 20, no. 1, pp. 63–76, 2019.
- [49] H. Jia, W. Zhang, R. Zheng, S. Wang, X. Leng, and N. Cao, "Ensemble mutation slime mould algorithm with restart mechanism for feature selection," *International Journal of Intelligent Systems*, vol. 37, no. 3, pp. 2335–2370, 2021.
- [50] K. K. Ghosh, R. Guha, S. K. Bera, N. Kumar, and R. Sarkar, "S-shaped versus V-shaped transfer functions for binary manta ray foraging optimization in feature selection problem," *Neural Computing & Applications*, vol. 33, no. 17, pp. 11027–11041, 2021.
- [51] J. Zhang, B. Gao, H. Chai, Z. Ma, and G. Yang, "Identification of DNA-binding proteins using multi-features fusion and binary firefly optimization algorithm," *BMC Bioinformatics*, vol. 17, no. 1, p. 323, 2016.
- [52] A. D. Li, B. Xue, and M. Zhang, "Multi-objective feature selection using hybridization of a genetic algorithm and direct multisearch for key quality characteristic selection," *Information Sciences*, vol. 523, pp. 245–265, 2020.
- [53] B. Abdollahzadeh, F. Soleimani Gharehchopogh, and S. Mirjalili, "Artificial gorilla troops optimizer: a new nature-inspired metaheuristic algorithm for global optimization problems," *International Journal of Intelligent Systems*, vol. 36, no. 10, pp. 5887–5958, 2021.