

Implementation of Features Selection Based on Dragonfly Optimization Algorithm

Nadia Mohammed Majeed¹

Nadia.20csp72@student.uomosul.edu.iq

Fawziya Mahmood Ramo¹

fawziyaramo@uomosul.edu.iq

¹Computer Science Department, University of Mosul. Mosul / Iraq

Abstract Nowadays increasing dimensionality of data produces several issues in machine learning. Therefore, it is needed to decrease the number of features by choosing just the most important ones and eliminating duplicate features, also reducing the number of features that are important to the model. For this purpose, many methodologies known as Feature Selection are applied. In this study, a feature selection approach is proposed based on Swarm Intelligence methods, which search for the best points in the search area to achieve optimization. In this paper, a wrapper feature selection technique based on the Dragonfly algorithm is proposed. The dragonfly optimization technique is used to find the optimal subset of features that could accurately classify breast cancer as benign or malignant. Many times, the fitness function is defined as classification accuracy. In this study, hard vote classes are employed as a model developed to evaluate feature subsets that have been chosen. It is used as an evaluation function (fitness function) to evaluate each dragonfly in the population. The proposed ensemble hard voting classifier utilizes a combination of five machine-learning algorithms to produce a binary classification for feature selection: Support Vector Machine (SVM), K-Nearest Neighbors (K-NN), Naive Bayes (NB), Decision Tree (DT), and Random Forest (RF). According to the results of the experiments, the voting ensemble classifier has the greatest accuracy value among the single classifiers. The proposed method showed that when training the subset features, the accuracy generated by the voting classifier is high at 98.24%, whereas the training of all features achieved an accuracy of 96.49%. The proposed approach makes use of the UCI repository's Wisconsin Diagnostic Breast Cancer (WDBC) Dataset. Which consists of 569 instances and 30 features.

Keywords: Dragonfly optimization algorithm, Feature Selection

1. Introduction

We are in the Big Data era. As the quantity of data we acquire grows, grows exponentially, extracting the essential data becomes increasingly essential. Data of high dimensions increases the search area and takes longer to process. When implementing model-based machine learning approaches. Furthermore, it might produce noise, which affects the model's building and reduces its efficacy. Researchers mostly employ two ways to overcome challenges resulting due to the data's high dimensionality: feature extraction and feature selection [1]. By collecting some attributes, feature extraction produces a feature space with a low dimension, while Feature selection eliminates duplicate features and reduces the number of features that are important to the model to a limited number. Feature Selection (FS) is an optimization issue that has a search area with all relevant attributes and the need to select the best ones. In an FS process, there are three essential steps: a technique for generating the next potential subset, a function for evaluating the subset under testing, and a procedure for selecting or abandoning the subset. In this paper, The Feature selection wrapping technique was proposed based on Swarm Intelligence. Among the many well-known SI algorithms, Dragonfly optimization was selected and a hard-voting classifier was employed to evaluate the feature subsets.[2].

2. Related Work

In previous years, several studies were conducted on the subject of ensemble learning, and multiple researches were presented regarding the applications of the Dragonfly algorithm to select features related to the output. Those studies that the researchers reached were summarized as follows:

Hazra et al. [3] proposed A comparison study was done comparing the various machine learning approaches employed, and their effective work helped them attain accuracy of 95.16 %, 95.53 %, and 95.91 %, respectively, for Naive Bayes, SVM, and Ensemble techniques.

Mafarja et al. [4] suggested the binary dragonfly method is employed as a wrapper-feature selection technique and as an evaluator, the K-NN classifier is used. The performance of the suggested technique is evaluated using eighteen UCI datasets. In regard to classification accuracy and the number of chosen attributes, the proposed method's results are compared to those of Particle Swarm Optimization (PSO) and Genetic Algorithms (GAs). The results demonstrate the Binary Dragonfly Algorithm's (BDA) potential to search the space of features that are most relevant for classification operations.

Yasen et al. [5] used the Dragonfly Algorithm (DA) to determine the weights of each link between Artificial Neural Network (ANN) neurons in a method called (ANN-DA). It uses real data as a case study and is used to forecast disease. Different measures have been used in an experimental comparison of both optimization techniques with other well-known classifiers such as Artificial Bee Colony (ABC), ANN-ABC. The findings demonstrate the superiority of ANN-DA over ANN-ABC and traditional ANN. Additionally; ANN-DA findings were more reliable among all datasets.

Feng et al. [6] suggested a combination technique (DA-SVM) for evaluating the short-term load forecasting of onshore oil field micro-grids in China's Bohai Sea, based on the hybridization of DA and SVM. In comparison to PSO-SVM, GA-SVM, and Back Propagation neural network (BPNN) methods, experimental findings showed that DA-SVM has a better global search capability and high predictability.

Nguyen et al. [7] proposed that The best breast cancer prediction model is ensemble voting, according to research. After selecting features, the data was used to test and train multiple classification models. Just four methods: AdaBoost, logistic regression (LR), SVM, and ensemble voting classifier, out of all the ones employed for prediction, perform better, with an accuracy of roughly 98% based on the findings.

Raza, Khalid. [8] Proposed an ensemble model that combines the output of three classification methods, including logistic regression, multilayer perceptron, naive Bayes, and to forecast heart disease, the majority voting method is used. The suggested ensemble technique obtained a classification accuracy of 88.88%, which is higher than any single classifier model.

Ling Li. [9] used a hybrid prediction model that included the SVM algorithm with an improved dragonfly algorithm to estimate short-term wind power. To increase the performance of the traditional dragonfly method, an adaptive learning factor and a differential evolution technique are proposed. The improved dragonfly method is used to choose the best SVM settings. On real datasets collected from France's La Haute Borne wind farm, the efficiency of the suggested model has been proven as compared with previous methods like back propagation neural networks and Gaussian process regression in terms of prediction accuracy.

MurtiRawat et al.[10] Proposed Various Machine Learning approaches, such as Logistic Regression (LR), (KNN), and voting classifier with Principal Component Analysis (PCA), have been proposed to aid in the detection of breast cancer. The models were trained and evaluated using data from the WDBC,. The data was pre-processed before being used to extract features from the data set utilizing Principal Component Analysis (PCA). The suggested technique had a classification accuracy of 98.60 % utilizing K-NN and 97.90 % utilizing Logistic Regression, whereas the voting classifier had the best accuracy of 99.30 %.

Assiri et al.[11] Proposed three classifiers for ensemble voting classifier: multilayer perceptron network, SVM learning with stochastic gradient descent optimization, and basic logistic regression learning. The majority-based voting mechanism was utilized for hard voting. When compared to the individual

method, the hard voting (majority-based voting) mechanism performs better with 99.42 % for the Wisconsin Breast Cancer Dataset WBCD.

Vinmalar, F. Leena, and A. Kumar Kombaiya. [12].proposed an Improved Dragonfly Algorithm (IDA) was utilized to decrease dataset's dimensionality of The lung cancer gene expression. To pick an effective subset of features, IDA employs the wrapper feature selection technique. For IDA feature selection and recognition of Lung types of tumours, the Random Stumors (RS), (ANN), and Sequential Minimal Optimization (SMO) classifiers were used. The fitness value of Dragonflies in each iteration is determined by the classifier's accuracy for a selected set of features of Dragonflies in the training dataset. Finally, the experimental study showed the IDA's accuracy, precision, recall, and F- measure effectiveness.

3. Research Methods

There are two kinds of meta-heuristic approaches: Evolutionary Algorithms (EA) and Swarm Intelligence Algorithms (SI). [13].

3.1. Swarm Intelligence

Swarm intelligence techniques are a collection of artificial intelligence methods Inspired it is behavior by nature concerned with the emergent collective of several cooperating agents adhere to a simple set of rules. Whereas each agent may be regarded as not intelligent, the entire system of several agents may exhibit self-organization behavior, allowing it to function as collective intelligence. Following are the main characteristics of Si-based algorithms and the following is a summary:

- 1- Agents collaborate and exchange information.
- 2- Each agent has self-management and autonomy.
- 3- Adaptability and the ability to react quickly to changes in the environment.
- 4- It is simple to parallelize for real-time issues

Each agents navigate in the search area and collaborate in SI algorithms, and represents a potential solution, which is assessed by the evaluation function. Agents improve with each iteration, finally finding a satisfactory solution. The fitness function's essential element is the machine-learning method's accuracy. The condition for stopping in algorithm optimization is usually a specified number of iterations or a result that is believed to be enough. As a result, even though neither of these conditions was met, one repetition of the main step is carried out. It involves moving agents, modifying their characteristics, and updating the better solution (s). The global better solution is eventually returned Swarm optimization techniques, like other search algorithms, require an appropriate balance of exploration and exploitation [2].

Each SI algorithm should follow a set of basic steps. As a result, as seen in (Figure 1), the SI framework is as follows:

- A. Initialize population.
- B. Establish a stop condition.
- C. Evaluate your fitness level.
- D. move and update agents.
- E. Provide the best global solution

The settings of the parameters of the algorithm must be determined before the initialization step.

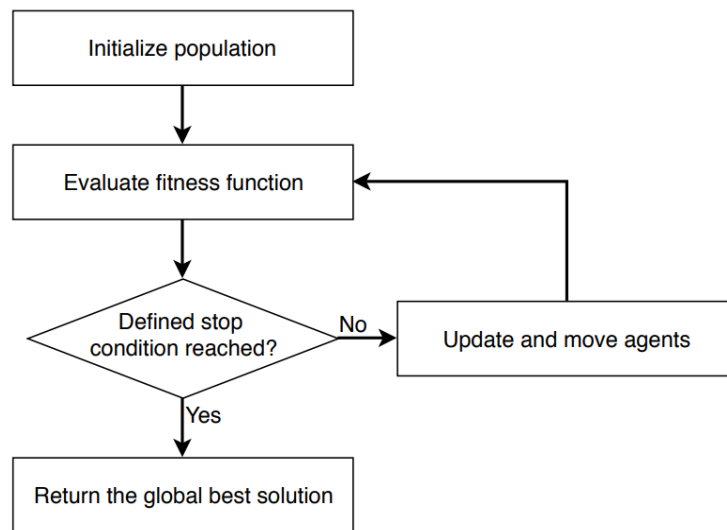


Figure 1: shows the framework for swarm intelligence

The best search agent is the output of a SI algorithm.[1]

3.2 Dragonfly Optimization Algorithm

Is a swarm-based advanced technique that was inspired by the static and dynamic grouping behaviours of dragonflies in nature. The technique is mathematically modeled by simulating the behaviours of dragonflies seeking prey [14]. By mimicking the hunting behaviour of dragonflies, the algorithm basically searches for the global optimal solution to the optimization issue. The life routines of dragonflies, such as seeking food, avoiding enemies, and determining flying paths are taken into account throughout the modelling.[15] Dragonflies have two major goals: hunting and migration. The first is known as a feeding (static) swarm, whereas the second is known as a migrating (dynamic) swarm.[16] In a dynamic swarm, A vast number of dragonfly groups move in common direction over lengthy distances in search of the best environment, but in a static swarm, each cluster consists of a few numbers of dragonflies.. a narrow range back and forth in search of other flying prey. The exploitation and exploration stages of meta-heuristics algorithm optimization can be compared to the migrating and feeding activities of dragonflies. In the exploitation phase, dragonflies congregate in big numbers and fly in a single direction, creating a dynamic swarm. In a static swarm, dragonflies in small groups fly back and forth. in a limited range to discover other flying prey, which is useful for search agent exploration. Figure (2,3) illustrates Mirjalili's dynamic and static dragonfly groupings.

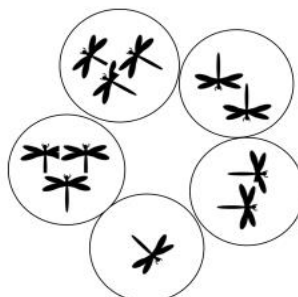


Fig 2. Static Swarm[17], [18]



Fig 3. Dynamic Swarm[17], [18]

In the insect swarms, separation, alignment, and cohesion are three basic characteristics. The degree of separation denotes an individual's static collision avoidance from other agents in the neighbourhood, the degree of alignment denotes an individual's velocity matching that of other agents in the neighbourhood, and the level of cohesiveness denotes an individual's tendency to gravitate toward the neighbourhood's center of mass. In DA, every swarm adheres to the survival principle, and each agent has two distinct behaviours: searching for food and attempting to avoid nearby enemies. Dragonflies use the following five behaviours to position themselves[14] :

1-Separation. Is the process that keeps the search agents in the neighbourhood apart. The separation behaviour is mathematically modeled in Equation (1):

$$S_i = - \sum_{j=1}^N (X_i - X_j)$$

(1)

Where S_i : individual's separation, X_i : individual's location, X_j : neighbouring agent's location, and N : The number of neighbourhoods.

2- Alignment. Refers to how one search agent's velocity compares to the other search agents' velocity in the area. Eq. (2) shows the mathematical formulation of the alignment behaviour:

$$A_i = \frac{\sum_{j=1}^N V_j}{N}$$

(2)

Where A_i : individual's alignment, V_j : surrounding individual's velocity, and N : The number of neighbourhoods.

3- Cohesion. Describes how individuals move from their neighbourhood to the center of mass. Individuals tend to fly towards the nearest center of mass. Eq. (3) shows the mathematical description of the Cohesion behaviour:

$$C_i = \frac{\sum_{j=1}^N X_j}{N} - X_i$$

(3)

Where C_i : individual's cohesion, X_i : individual's position, N : denotes the number of neighbouring individuals, and X_j : the position of the neighbour individual.

4- Attraction. The following formula calculates the attraction to the food source:

$$F_i = X^+ - X_i$$

(4)

Where F_i : individual's food source, X_i : individual's location, and X^+ : the food source's location.

5- Distraction. The following is how an enemy's distraction is calculated:

$$E_i = X^- + X_i$$

(5)

Where E_i : the enemy of the individual's position, X_i : individual's location, and X^- : the natural enemy's location.

Figure (4,5) illustrates the previous five swarming behaviours in dragonfly positioning movement.

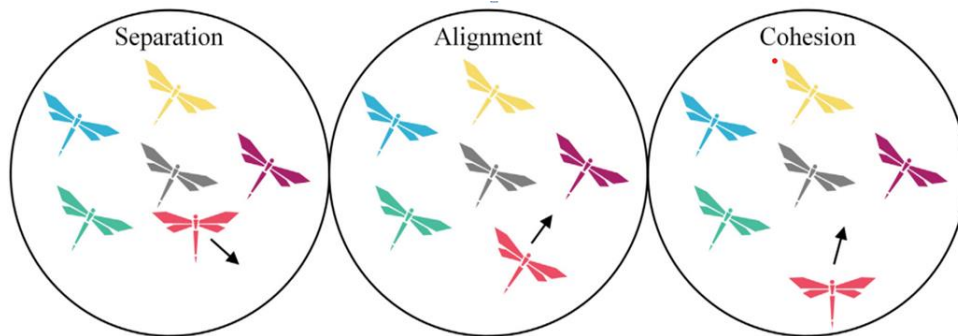


Fig. 4. Nature's three elements of swarming are separation, alignment, and cohesiveness.

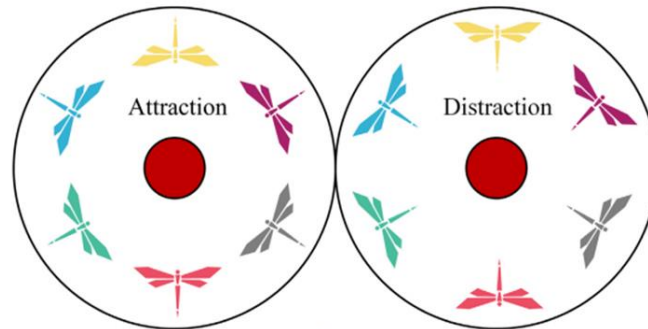


Fig. 5. Attraction and distraction are two additional natural swarming concepts discussed in DA.[19]

The position of dragonflies in a search area is updated using two vectors and to emulate their movements: step vector (ΔX) and position vector (X). The step vector indicates the travel direction of dragonflies can be described as follows:

$$\Delta X_i^{t+1} = (sS_i + aA_i + cC_i + fF_i + eE_i) + \omega\Delta X_i^t \quad (6)$$

Where S : separation weight, S_i : the i -th individual's separation, A : alignment weight, A_i : alignment of the i -th individual, and c : weight cohesion. C_i : the i -th individual's cohesiveness, f : food factor, F_i : the i -th individual's food source, e : enemy factor, E_i : enemy of the i -th individual, w : weight of inertia, and t : the number of iterations count.[14], [20].

Following the completion of the step vector computation, the update of the position vector is as follows

$$X_i^{t+1} = X_i^t + \Delta X_i^{t+1} \quad (7)$$

The current iteration number is denoted by the t .

The neighbours of each dragonfly are regarded by supposing a radius around each of them. The radius of the neighbourhoods is increased proportionately to the iteration counter to move from the discovery stage to the exploitation stage. As a result, static swarms are transformed into dynamic swarms. In the final step of optimization, all of the dragonflies will unite into a single dynamic swarm that will lead to the global optimal solution. When artificial dragonflies have no neighbours, they employ the Lévy flying mechanism [7] to navigate around the search area. A random walk is utilized to create a random position for dragonflies that do not have any neighbours. The position update formula utilized in this case is:

$$X_i^{t+1} = X_i^t + Levy(dim) \times X_i^t \quad (8)$$

Where d is the dimension of the position vectors and t is the current iteration number. Each dragonfly's step vector and position vectors are updated in each loop until the final criteria is satisfied reached. The following is a description of the Lévy function.

$$Levy(dim) = 0.01 \times \frac{r_1 \times \sigma}{|r_2|^{\frac{1}{\beta}}} \quad (9)$$

r_1 and r_2 are random values in the range [0,1], σ is a constant, and: [14]

$$\sigma = \left\{ \frac{\Gamma(1+\beta) \times \sin\left(\frac{\pi\beta}{2}\right)}{\Gamma\left(\frac{1+\beta}{2}\right) \times \beta \times 2^{(\beta-1)/2}} \right\}^{\frac{1}{\beta}} \quad (10)$$

3.3. Feature Selection

The most essential stage of feature selection and relevant features to be employed in model training, where feature selection goes through a filtering stage to exclude several features which are not necessary[21]. In classification, feature selection is a fundamental pre-processing procedure. The main benefits of the feature selection process are enhancing the model's capabilities and reducing the computing cost. Employing raw attributes may result in an ineffective output. While optimal attributes selection will play an important role in effective forecasting. As a result, various techniques utilize a variety of attributes selection strategies to prepare relevant data for the model training.[22]. To select the optimum feature subset that may be utilized to separate data records for a dataset with N characteristics, 2N subsets must be tested.[20] Hence, finding all feasible subsets exhaustively becomes impracticable and computationally costly for high-dimensional data. Search algorithms that may be employed to increase the efficiency of the features selection method include backwards, forwards, random, and heuristic searching Recently many FS approaches have been effectively utilized with metaheuristics algorithms [23]. For an explanation of a wide region solution space, the metaheuristic optimization algorithm is utilized[24]. Swarm Intelligence (SI) is a set of metaheuristic methods inspired by natural behaviours such as ant, bee, bird, and dragonfly cooperative behaviour, etc.

There are three types of FS algorithms[25]:

1. Wrapper models
2. Filter models
3. Embedded Models

A classifier is utilized and trained in the wrapper technique to evaluate a set of prominent features. A Particular learning model is utilized in the wrapper model to assess a subset of attributes in their search operations to select the set of features with better precision in classification. Many wrapper approaches use iterative search methods, in which the population of solutions is directed by each iteration of the learning model. to the optimal solution.[13].

4. Result and discussion

This research suggest a novel wrapper feature selection method using the dragonfly algorithm. The suggested FS approach's major objective was to determine the minimum reduction to achieve greater accuracy than using all of the dataset's features. The suggested technique is evaluated using a collection of well-known FS datasets from the UCI data repository. This dataset contains 569 instances, each of which has 30 features obtained from nucleus pictures and is separated into 30 % and 70 % testing and training, respectively. The 17 specific features obtained are compared with all the 30 features before the selection process, which are evaluated using a hard voting classifier the result shows that the accuracy of selected features is greater than the accuracy with all features as shown in the table below:

Tab. 1. All features and selected features obtained from prediction metrics

	No. of features	Accuracy	Precision	Recall	F1_Score
All Features	30	96.49122	95.53571	99.07407	97.27272
Selected Features	17	98.24561	97.29729	100.0	98.63013

5. Conclusion and future work

The purpose of this research is to classify the two categories of breast cancer tumours: benign and malignant tumours. In this research, a wrapper feature selection has suggested a strategy based on the Dragonfly optimization technique, which is used to create a subset of features that could accurately classify breast cancer. In this study, the hard vote classifier is used as a fitness function to evaluate feature subsets that have been chosen by the dragonfly optimization algorithm. The experiment's result showed that the accuracy produced by the voting classifier for selected features was higher than the accuracy of all features. In future work, it is worth experimenting with various evaluation functions to observe how DA behaves.

References

- [1] L. Brezočnik, I. Fister, and V. Podgorelec, "Swarm intelligence algorithms for feature selection: A review," *Appl. Sci.*, vol. 8, no. 9, 2018, doi: 10.3390/app8091521.
- [2] G. Kicska and A. Kiss, "Comparing swarm intelligence algorithms for dimension reduction in machine learning," *Big Data Cogn. Comput.*, vol. 5, no. 3, 2021, doi: 10.3390/bdcc5030036.
- [3] A. Hazra, S. Kumar, and A. Gupta, "Study and Analysis of Breast Cancer Cell Detection using Naïve Bayes, SVM and Ensemble Algorithms," *Int. J. Comput. Appl.*, vol. 145, no. 2, pp. 39–45, 2016, doi: 10.5120/ijca2016910595.
- [4] M. M. Mafarja, D. Eleyan, I. Jaber, A. Hammouri, and S. Mirjalili, "Binary Dragonfly Algorithm for Feature Selection," *Proc. - 2017 Int. Conf. New Trends Comput. Sci. ICTCS 2017*, vol. 2018-Janua, pp. 12–17, 2017, doi: 10.1109/ICTCS.2017.43.
- [5] T. Xie, J. Yao, and Z. Zhou, "DA-based parameter optimization of combined kernel support vector machine for cancer diagnosis," *Processes*, vol. 7, no. 5, 2019, doi: 10.3390/pr7050263.
- [6] Y. Feng, P. Zhang, M. Yang, Q. Li, and A. Zhang, "Short term load forecasting of offshore oil field microgrids based on DA-SVM," *Energy Procedia*, vol. 158, pp. 2448–2455, 2019, doi: 10.1016/j.egypro.2019.01.318.
- [7] Q. H. Nguyen et al., "Breast Cancer Prediction using Feature Selection and Ensemble Voting," *Proc. 2019 Int. Conf. Syst. Sci. Eng. ICSSE 2019*, pp. 250–254, 2019, doi: 10.1109/ICSSE.2019.8823106.
- [8] K. Raza, *Improving the prediction accuracy of heart disease with ensemble learning and majority voting rule*. Elsevier Inc., 2019.
- [9] L. L. Li, X. Zhao, M. L. Tseng, and R. R. Tan, "Short-term wind power forecasting based on support vector machine with improved dragonfly algorithm," *J. Clean. Prod.*, vol. 242, p. 118447, 2020, doi: 10.1016/j.jclepro.2019.118447.
- [10] R. Murtirawat, S. Panchal, V. K. Singh, and Y. Panchal, "Breast Cancer Detection Using K-Nearest Neighbors, Logistic Regression and Ensemble Learning," *Proc. Int. Conf. Electron. Sustain. Commun. Syst. ICESC 2020*, no. Icesc, pp. 534–540, 2020, doi: 10.1109/ICESC48915.2020.9155783.
- [11] A. S. Assiri, S. Nazir, and S. A. Velastin, "Breast Tumor Classification Using an Ensemble Machine Learning Method," *J. Imaging*, vol. 6, no. 6, 2020, doi: 10.3390/JIMAGING6060039.
- [12] F. L. vinmalar* and D. A. K. Kombaiya, "An Improved Dragonfly Optimization Algorithm based Feature Selection in High Dimensional Gene Expression Analysis for Lung Cancer Recognition," *Int. J. Innov. Technol. Explor. Eng.*, vol. 9, no. 8, pp. 896–908, 2020, doi: 10.35940/ijitee.h6302.069820.
- [13] M. Rostami, K. Berahmand, E. Nasiri, and S. Forouzande, "Review of swarm intelligence-based feature selection methods," *Eng. Appl. Artif. Intell.*, vol. 100, no. January, p. 104210, 2021, doi: 10.1016/j.engappai.2021.104210.
- [14] L. Wang, R. Shi, and J. Dong, "A hybridization of dragonfly algorithm optimization and angle modulation mechanism for 0-1 knapsack problems," *Entropy*, vol. 23, no. 5, pp. 1–24, 2021,

- doi: 10.3390/e23050598.
- [15] Y. Yue et al., “A Data Collection Method for Mobile Wireless Sensor Networks Based on Improved Dragonfly Algorithm,” *Comput. Intell. Neurosci.*, vol. 2022, pp. 1–16, 2022, doi: 10.1155/2022/4735687.
- [16] N. Devarakonda, S. Anandarao, R. Kamarajugadda, and Y. Wang, “UNIQUE DRAGONFLY OPTIMIZATION ALGORITHM FOR,” 2019.
- [17] M. Alshinwan et al., “Dragonfly algorithm: a comprehensive survey of its results, variants, and applications,” *Multimed. Tools Appl.*, pp. 14979–15016, 2021, doi: 10.1007/s11042-020-10255-3.
- [18] C. M. Rahman and T. A. Rashid, “Dragonfly algorithm and its applications in applied science survey,” *Comput. Intell. Neurosci.*, vol. 2019, 2019, doi: 10.1155/2019/9293617.
- [19] J. Too and S. Mirjalili, “A Hyper Learning Binary Dragonfly Algorithm for Feature Selection: A COVID-19 Case Study,” *Knowledge-Based Syst.*, vol. 212, p. 106553, 2021, doi: 10.1016/j.knsys.2020.106553.
- [20] A. I. Hammouri, M. Mafarja, M. A. Al-Betar, M. A. Awadallah, and I. Abu-Doush, “An improved Dragonfly Algorithm for feature selection,” *Knowledge-Based Syst.*, vol. 203, p. 106131, 2020, doi: 10.1016/j.knsys.2020.106131.
- [21] A. Abdulmunim Abdulmajeed Althanoon and Y. S. Younis, “Supporting Classification of Software Requirements system Using Intelligent Technologies Algorithms,” *Tech. Rom. J. Appl. Sci. Technol.*, vol. 3, no. 11, pp. 32–39, 2021, doi: 10.47577/technium.v3i11.5417.
- [22] W. Shafqat, S. Malik, K. T. Lee, and D. H. Kim, “Pso based optimized ensemble learning and feature selection approach for efficient energy forecast,” *Electron.*, vol. 10, no. 18, 2021, doi: 10.3390/electronics10182188.
- [23] I. Zelinka, “A survey on evolutionary algorithms dynamics and its complexity - Mutual relations, past, present and future,” *Swarm Evol. Comput.*, vol. 25, pp. 2–14, 2015, doi: 10.1016/j.swevo.2015.06.002.
- [24] H. M. Osman, R. S. Alsawaf, and A. Y. Hammo, “Survey of using grasshopper algorithm,” *Tech. Rom. J. Appl. Sci. Technol.*, vol. 4, no. 3, pp. 37–44, 2022, doi: 10.47577/technium.v4i3.6344.
- [25] C. C. Aggarwal, X. Kong, Q. Gu, J. Han, and P. S. Yu, “Active learning: A survey,” *Data Classif. Algorithms Appl.*, pp. 571–605, 2014, doi: 10.1201/b17320.