

# **Classification and Feature Selection Method for Medical Datasets by BGEO -TVFL(Binary golden eagle optimization-Time Varying Flight length) and KNN(k-nearest neighbour)**

*A Project Report submitted in the partial fulfillment of the Requirements for the award of the degree*

## **BACHELOR OF TECHNOLOGY IN COMPUTER SCIENCE AND ENGINEERING**

Submitted by

**B. Lakshmi Prasanna (21471A0513)**  
**P. Swetha Lakshmi (21471A0543)**  
**M. Sahithi (21471A0535)**

Under the esteemed guidance of  
**Dr. Rama Krishna Eluri, M.Tech., Ph.D.**  
Assistant Professor



## **DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

**NARASARAOPETA ENGINEERING COLLEGE: NARASAROPET**

**(AUTONOMOUS)**

Accredited by NAAC with A+ Grade and NBA under Tyre -1

NIRF rank in the band of 201-300 and an ISO 9001:2015 Certified

Approved by AICTE, New Delhi, Permanently Affiliated to JNTUK, Kakinada  
KOTAPPAKONDAROAD, YALAMANDAVILLAGE, NARASARAOPET-  
522601

2024-2025

NARASARAOPETAENGINEERING COLLEGE  
(AUTONOMOUS)

DEPARTMENT OF COMPUTERSCIENCEAND  
ENGINEERING



CERTIFICATE

This is to certify that the project that is entitled with the name Classification and Feature Selection Method for Medical Datasets by BGEO TVFL (Binary golden eagle optimization-Time Varying Flight length) and KNN(k-nearest neighbour is a bonafide work done by the team **Bommu Lakshmi Prasanna (21471A0513)**, **Paruchuri Swetha Lakshmi(21471A0543)**, **Mellacheruvu Sahithi(21471A0535)** in partial fulfillment of the requirements for the award of the degree of BACHELOR OF TECHNOLOGY in the Department of COMPUTER SCIENCE AND ENGINEERING during 2024-2025.

PROJECT GUIDE

**Dr. Rama Krishna Eluri, M.Tech., Ph.D.**  
Assistant Professor

PROJECT CO-ORDINATOR

**D.Venkata Reddy, B.Tech. M.Tech. (Ph.D.)**  
Assistant Professor

HEAD OF THE DEPARTMENT

**Dr. S. N. Tirumala Rao, M.Tech., Ph.D.,**  
Professor & HOD

EXTERNALEXAMINER

## **DECLARATION**

We declare that this project work titled “CLASSIFICATION AND FEATURE SELECTION METHOD FOR MEDICAL DATASETS BY BGEO TVFL (BINARY GOLDEN EAGLE OPTIMIZATION-TIME VARYING FLIGHT LENGTH) AND KNN (K-NEAREST NEIGHBOUR)” is composed by ourselves that the work contain here is our own except where explicitly stated otherwise in the text and that this work has been submitted for any other degree or professional qualification except as specified.

B. Lakshmi Prasanna (21471A0513)  
P. Swetha Lakshmi (21471A0543)  
M. Sahithi (21471A0535)

## **ACKNOWLEDGEMENT**

We wish to express my thanks to various personalities who are responsible for the completion of the project. We are extremely thankful to our beloved chairman **sri M. V. Koteswara Rao, B.Sc.**, who took keen interest in us in every effort throughout this course. We owe our sincere gratitude to our beloved principal **Dr. S. Venkateswarlu, Ph.D.**, for showing his kind attention and valuable guidance throughout the course.

We express our deep-felt gratitude towards **Dr. S. N. Tirumala Rao, M.Tech., Ph.D.**, HOD of CSE department and also to our guide **Dr. Rama Krishna Eluri, M.Tech., Ph.D.**, of CSE department whose valuable guidance and unstinting encouragement enable us to accomplish our project successfully in time.

We extend our sincere thanks to **D. Venkat Reddy B.Tech, M.Tech.,(Ph.D.)**, Assistant professor & Project coordinator of the project for extending his encouragement. Their profound knowledge and willingness have been a constant source of inspiration for us throughout this project work.

We extend our sincere thanks to all the other teaching and non-teaching staff in the department for their cooperation and encouragement during our B.Tech degree. We have no words to acknowledge the warm affection, constant inspiration and encouragement that we received from our parents.

We affectionately acknowledge the encouragement received from our friends and those who were involved in giving valuable suggestions clarified our doubts, which really helped us in successfully completing our project.

By

**B.Lakshmi Prasanna (21471A0513)**  
**P.Swetha Lakshmi (21471A0543)**  
**M.Sahithi (21471A0535)**



## **INSTITUTE VISION AND MISSION**

### **INSTITUTION VISION**

To emerge as a Centre of excellence in technical education with a blend of effective student centric teaching learning practices as well as research for the transformation of lives and community.

### **INSTITUTION MISSION**

M1: Provide the best class infra-structure to explore the field of engineering and research

M2: Build a passionate and a determined team of faculty with student centric teaching, imbibing experiential, innovative skills

M3: Imbibe lifelong learning skills, entrepreneurial skills and ethical values in students for addressing societal problems



## **DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

### **VISION OF THE DEPARTMENT**

To become a centre of excellence in nurturing the quality Computer Science & Engineering professionals embedded with software knowledge, aptitude for research and ethical values to cater to the needs of industry and society.

### **MISSION OF THE DEPARTMENT**

The department of Computer Science and Engineering is committed to

**M1:** Mould the students to become Software Professionals, Researchers and Entrepreneurs by providing advanced laboratories.

**M2:** Impart high quality professional training to get expertise in modern software tools and technologies to cater to the real time requirements of the Industry.

**M3:** Inculcate team work and lifelong learning among students with a sense of societal and ethical responsibilities.



## **Program Specific Outcomes (PSO's)**

**PSO1:** Apply mathematical and scientific skills in numerous areas of Computer Science and Engineering to design and develop software-based systems.

**PSO2:** Acquaint module knowledge on emerging trends of the modern era in Computer Science and Engineering

**PSO3:** Promote novel applications that meet the needs of entrepreneur, environmental and social issues.



## **Program Educational Objectives (PEO's)**

The graduates of the programme are able to:

**PEO1:** Apply the knowledge of Mathematics, Science and Engineering fundamentals to identify and solve Computer Science and Engineering problems.

**PEO2:** Use various software tools and technologies to solve problems related to academia, industry and society.

**PEO3:** Work with ethical and moral values in the multi-disciplinary teams and can communicate effectively among team members with continuous learning.

**PEO4:** Pursue higher studies and develop their career in software industry.

## Program Outcomes

**PO1: Engineering knowledge:** Apply the knowledge of mathematics, science, engineering fundamentals, and an engineering specialization to the solution of complex engineering problems.

**PO2: Problem analysis:** Identify, formulate, research literature, and analyse complex engineering problems reaching substantiated conclusions using first principles of mathematics, natural sciences, and engineering sciences.

**PO3: Design/development of solutions:** Design solutions for complex engineering problems and design system components or processes that meet the specified needs with appropriate consideration for the public health and safety, and the cultural, societal, and environmental considerations.

**PO4: Conduct investigations of complex problems:** Use research-based knowledge and research methods including design of experiments, analysis and interpretation of data, and synthesis of the information to provide valid conclusions.

**PO5: Modern tool usage:** Create, select, and apply appropriate techniques, resources, and modern engineering and IT tools including prediction and modeling to complex engineering activities with an understanding of the limitations.

**PO6: The engineer and society:** Apply reasoning informed by the contextual knowledge to assess societal, health, safety, legal and cultural issues and the consequent responsibilities relevant to the professional engineering practice.

**PO7: Environment and sustainability:** Understand the impact of the professional engineering solutions in societal and environmental contexts, and demonstrate the knowledge of, and need for sustainable development

**PO8: Ethics:** Apply ethical principles and commit to professional ethics and responsibilities and norms of the engineering practice.

**PO9: Individual and team work:** Function effectively as an individual, and as a member or leader in diverse teams, and in multidisciplinary settings.

**PO10: Communication:** Communicate effectively on complex engineering activities with the engineering community and with society at large, such as, being able to comprehend and write effective reports and design documentation, make effective presentations, and give and receive clear instructions.

**PO11: Project management and finance:** Demonstrate knowledge and understanding of the engineering and management principles and apply these to one's own work, as a member and leader in a team, to manage projects and in multidisciplinary environments.

**PO12: Life-long learning:** Recognize the need for, and have the preparation and ability to engage in independent and life-long learning in the broadest context of technological change.



### **Project Course Outcomes (CO'S):**

**CO421.1:** Analyse the System of Examinations and identify the problem.

**CO421.2:** Identify and classify the requirements.

**CO421.3:** Review the Related Literature

**CO421.4:** Design and Modularize the project

**CO421.5:** Construct, Integrate, Test and Implement the Project.

**CO421.6:** Prepare the project Documentation and present the Report using appropriate method.

### **Course Outcomes – Program Outcomes mapping**

	IO 1	PO2	PO 3	PO 4	IO 5	PO 6	PO 7	PO 8	IO 9	PO10	PO11	PO12	PSO1	PSO2	PSO3
<b>C421.1</b>		✓											✓		
<b>C421.2</b>	✓		✓		✓								✓		
<b>C421.3</b>				✓		✓			✓				✓		
<b>C421.4</b>			✓			✓		✓					✓	✓	
<b>C421.5</b>					✓	✓		✓	✓	✓	✓	✓	✓	✓	✓
<b>C421.6</b>									✓	✓	✓		✓	✓	

## Course Outcomes – Program Outcome correlation

	PO1	PO2	PO3	PO4	PO5	PO6	PO7	PO8	PO9	PO10	PO11	PO12	PSO1	PSO2	PSO3
<b>C421.</b>	2	3											2		
<b>C421.</b>			2		3								2		
<b>C421.</b>				2		2	3	3					2		
<b>C421.</b>				2		1	1	2					3	2	
<b>C421.</b>					3	3	3	2	3	2	2	1	3	2	1
<b>C421.</b>									3	2	1		2	3	

**Note: The values in the above table represent the level of correlation between CO's and PO's:**

1. Lowlevel
2. Medium level
3. High level

## Project mapping with various courses of Curriculum with Attained PO's:

Name of the course from which principles are applied in this project	Description of the device	Attained PO
C2204.2, C22L3.2	Gathering requirements and defining the problem, planning to develop a model for breast cancer detection using BGEO-TVFL with KNN.	PO1, PO3
CC421.1,C2204.3, C22L3.2	Critical analysis of each requirement and identification of an optimized feature selection and classification model.	PO2, PO3
CC421.2,C2204.2, C22L3.3	Logical design is done by using the unified modelling language which involves individual team work	PO3,PO5, PO9
CC421.3,C2204.3, C22L3.2	Each and every module is tested, integrated, and evaluated in our projects	PO1, PO5
CC421.4,C2204.4, C22L3.2	Documentation is done by all our four members in the form of a group	PO10
CC421.5,C2204.2, C22L3.3	Each and every phase of the work in group is presented periodically	PO10, PO11
C2202.2, C2203.3, C1206.3,C3204.3, C4110.2	Implementation of the <b>BGEO-TVFL and KNN-based medical data classification system</b> , enabling healthcare professionals to efficiently analyze medical datasets. The system optimizes <b>feature selection</b> to improve classification accuracy while reducing computational complexity. Future enhancements may extend the model's capability to classify <b>other medical conditions</b> ,	PO4, PO7
C32SC4.3	The physical design includes website to check whether an image is real or fake	PO5, PO6

## ABSTRACT

Classification accuracy and feature selection are important steps in medical data analysis to identify suitable features that can improve the performance of machine learning models. In this study, we present a new method called BGEO TVFL (Binary Golden Eagle Optimization-Time Varying Flight Length) algorithm together is proposed for its K-nearest neighbour (KNN) algorithm. The TVFL algorithm feature of BGEO is used for optimal subset selection, while the KNN algorithm is used for classification. The proposed method is tested in several projects. The experimental results show that the proposed method outperforms other existing methods such as BWOA, BGWO, ACO, and ABC in terms of accuracy and selected features. Our results show that BGEO TVFL outperforms other algorithms in terms of accuracy and feature selection, achieving higher classification accuracy and selecting fewer features compared to how BGEO TVFL performs well as a good optimization algorithm for feature selection and classification in medical data sets.

**KEYWORDS:** Feature Selection, BGEO (Binary Golden Eagle Optimization), Classification, TVFL (Time Varying Flight Length), Wrapper Method

## **INDEX**

<b>S.NO.</b>	<b>CONTENT</b>	<b>PAGE NO</b>
1	INTRODUCTION	1-7
2	LITERATURE SURVEY	8-12
3	SYSTEM ANALYSIS	13
	3.1 EXISTING SYSTEM	14
	3.2 DISADVANTAGES OF EXISTING SYSTEM	16
	3.3 PROPOSED SYSTEM	16
	3.4 FEASIBILITY STUDY	20
	3.5 USING COCOMO MODEL	21
4	SYSTEM REQUIREMENTS	23
	4.1 SOFTWARE REQUIREMENTS	23
	4.2 HARDWARE REQUIREMENTS	24
	4.3 SOFTWARE DESCRIPTION	25
5	SYSTEM DESIGN	26
	5.1 SYSTEM ARCHITECTURE	26
	5.2 MODULES	29
	5.3 UML DIAGRAMS	33
6	IMPLEMENTATION	37
	6.1 MODEL IMPLEMENTATION	37
	6.2 CODING	51
7	TESTING	67
	7.1 TYPES OF TESTING	67
	7.2 INTEGRATION TESTING	69
8	RESULT ANALYSIS	71-78
9	USER INTERFACE	79-82
10	CONCLUSION	83
11	FUTURE SCOPE	84
12	REFERENCES	85-86

<b>S.NO.</b>	<b>LIST OF FIGURES</b>	<b>PAGE NO</b>
1.	Fig 1 General Process of FS	03
2.	Fig 2 Flow chart for BGEO-TVFL	19
3.	Fig3 Use case Diagram	33
4.	Fig4 Class Diagram	34
5.	Fig5 Sequence Diagram	35
6.	Fig6 Activity Diagram	36
7.	Fig 7 Accuracy graphs for all datasets	75
8.	Fig 7.1 Accuracy graph for breast_cancer.csv	75
9.	Fig 7.2 Accuracy graph for heart.csv	75
10.	Fig 7.3 Accuracy graph for Iris.csv	75
11.	Fig 7.4 Accuracy graph for diabetes.csv	75
12.	Fig 7.5 Accuracy graph for thyroid_csv	76
13.	Fig 8 Classification graphs for all datasets	77
14.	Fig 8.1 Classification graphs for breast_cancer	77
15.	Fig 8.2 Classification graphs for heart.csv	77
16.	Fig 8.3 Classification graphs for Iris.csv	77
17.	Fig 8.4 Classification graphs for diabetes.csv	77
18.	Fig 8.5 Classification graphs for thyroid_csv	77
19.	Fig 9 Radar Graph for all algorithms with all datasets	78
20.	Fig 10 Home Page	79
21.	Fig 11 About our project	79
22.	Fig 12 Upload your Datasets	80
23.	Fig 13 Result Page	80
24.	Fig 14 Flowchart	81
25.	Fig 15 Model Evaluation Metrics	81
26.	Fig 16 Model Evaluation Metrics	82

<b>S.NO.</b>	<b>LIST OF TABLES</b>	<b>PAGE NO</b>
1	Dataset Description	71
2	Outcomes of the proposed BGEO-TVFL	72
3	Comparison of classification accuracy with Proposed BGEO-TVFL and existing approaches	74
4	Comparison of selected features with Proposed BGEO-TVFL and existing approaches	76

# 1. INTRODUCTION

Feature selection is an important step in machine learning [8], especially in medical datasets [1] where the number of features may be too large and irrelevant features may lead to poor model performance. Feature selection identifies relevant informative features of large candidates and poor generalization. The selection problem has been extensively studied in a variety of fields including medicine [7], where the complexity of medical data sets requires effective and efficient selection strategies, which may not be effective in tight correlation considering the differences between the components of the medical database. To address this limitation, this study proposed a new selection method called binary Golden Eagle optimization-time-varying flight length (BGEO-TVFL), with time-varying capability of binary optimization methods can optimize the feature selection process[5] -By taking advantage of the K-Nearest Neighbour (KNN) algorithm as a classifier combining flight lengths, we show that BGEO-TVFL outperforms other state-of-the-art approaches optimize features such as binary whale optimization algorithm (BWOA)[6], binary grey wolf optimization (BGWO)[17], Ant Colony Optimization (ACO),[24] and Artificial Bee Colony (ABC)[23] in feature selection accuracy and computational efficiency. However, their performance should be evaluated against other optimization algorithms for feature selection. In this paper, we propose a new BGO called Binary Golden Eagle Optimization-Time Varying Flight Length (BGEO TVFL), which incorporates time-varying flight length to improve the search process We use BGEO TVFL compare for options in the medical lists BWOA, BGWO, . Let's do it with ACO, and the Artificial Bee Colony (ABC) algorithm. Our results show that BGEO TVFL outperforms other algorithms in terms of accuracy and feature selection.

## 1.1 Feature Selection

One of key steps in feature engineering is selection of feature, which involves choosing the most crucial characteristics to use in machine learning algorithms (Tang and Mahmoud, 2021). FS strategies are employed for minimizing the number of input variables by eliminating duplicate or unneeded features. Compressing collection of features to the one feature is a most significant to ML model (Das, Pradhan, Mishra, Pradhan and Patnaik, 2022). Each and every ML procedure is dependent on feature engineering, which basically entails the two processes, extraction of features and selecting features. Although the

Hybrid cluster Selection of features and extraction approaches may have the same end in mind, they are very different from one another. The main difference between the two is that feature extraction adds new features whereas feature selection just chooses a subset of the initial feature set. FS is a reducing the input variable method for a model by choosing relevant information in order to decrease overfitting in the model.

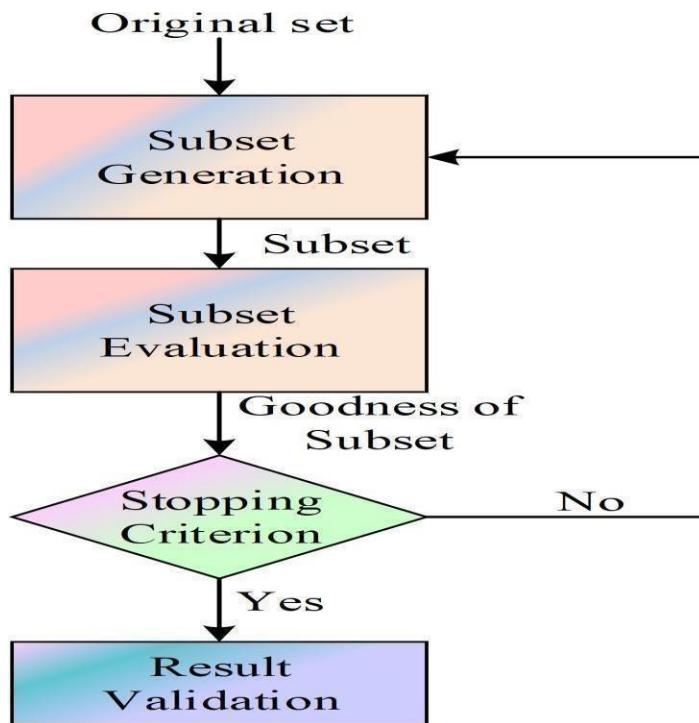
### **1.1.1 Need of FS**

It is crucial to comprehend the requirement for any approach before using it, as well as the need for selecting features [3]. For receiving better outcomes from ML, a high-quality and pre- processed input dataset is required (Nafis and Awang, 2021). For improving the system and aid in its learning, there is need of gathering a vast quantity of data. Frequently, the dataset consists of noisy data, worthless data, and a little quantity of useful data. The model's training process is also slowed down by the vast amount of data, and it is unable to perform well or make accurate predictions when faced with noise and unrelated input. As a result, it is essential to remove unwanted disturbances and irrelevant data from the dataset. Techniques for choosing features are used to accomplish this. The main benefits of completing FS beforehand rather than letting the ML model select which features are most important is discussed underneath (Amjad et al., 2021):

- Simpler models: Un useful models are those that are extremely complicated and challenging to comprehend. Additionally, simpler models are simpler to express.
- Fewer training cycles: a model requires less training time when the subset of features is more exact.
- Reduce variance: improve the precision of estimations that may be extracted for a specific simulation.
- Avoid the high-dimensionality curse: The volume of space increases so fast as size and quantity of features maximize that the quantity of data is limited.

### 1.1.2 Process of FS

In many cases, irrelevant and redundant characteristics or noise in the data may impede the feature selection process since they are not relevant and crucial in terms of the class idea, such as data analysis of microarray (Huang, 2021). When there are much less samples than features, ML becomes particularly difficult since the search space will be poorly filled. As a consequence, the model won't be able to accurately differentiate between important data and noise (Scatterer et al., 2022). There are two main methods for choosing features. The first is called Individual Evaluation, while the second is called Subset Evaluation.



**Fig. 1. General Process of FS**

Individual evaluation is the term used to rank the features (Khaire and Dhanalakshmi, 2022). In individual evaluation, a feature's weight is allocated in accordance with how relevant it is. Candidate feature subsets are created in subset evaluation using a search approach.

According to Fig. 1, there are four essential phases in the overall process for feature selection.

- Generation of Subset
- Estimation of Subset

- Criteria for Stopping

- Validation of Result

A heuristic search called subset generation selects a potential subset for consideration in the search space for each state. The nature of process of subset formation is determined by two fundamental considerations. The search direction is influenced by the search beginning point, which is first decided by the succeeding generation. Backward, forward, weighted, compound, and random methods may be taken into account when determining the search beginning locations at each stage (Akinola et al., 2022). Additional, the process of selecting features are handled by the search organization, which employs a specific plan such as exponential, sequential, or random searches. A newly generated subset must meet certain requirements in order to be accepted. Several criteria are proposed in literature to evaluate the level of quality of the subset candidates. Evaluation measures may be separated into dependent and independent criteria based on how much they depend on mining methods. Without using any mining methods, independent condition utilize the key traits of the data for training to measure a feature set's or feature's worth. Additionally, dependent criteria include programmed mining methods for feature selection [4] that choose features depending on how well they work when used on the traits that have been selected as a subset. Finally, in order to conclude the selection process, stop criteria must be created. The validation procedure is where the feature selection process ends. Even though it is not a essential stage in the feature selection (FS) process, the cogency of the FS technique must be proven by numerous tests and relationships with findings that have already been began or with those of challenging methods utilizing forged datasets, real-world datasets, or both.

A model is inferred from the interaction of the FS algorithm and the inductive learning approach. For FS, there are three general methods. The Filter Approach[24] first takes use of the common traits of training data without relying on the mining algorithm (Masood et al., 2023). The Wrapper Approach [24,26] also investigates the connection between choosing the best feature subset and relevance. It looks for an ideal feature subset that is tailored to the particular mining algorithm (Balakrishnan et al., 2022). Third, the Embedded Approach employs a particular learning algorithm to carry out feature selection during training.

## 1.2 Problem Definition

The problem of feature selection arises in many areas of data analysis where a large number of input variables must be considered to make predictions or decisions. Various feature selection algorithms have been developed, including wrapper methods that utilize optimization algorithms to identify issues and improve performance. However, optimization algorithms used in wrapper methods often have specific restrictions that must be tuned to improve performance, which is a difficult and time-consuming process. For example, PSO[17] and GA have four and six parameters that required to be tuned, respectively, and other algorithms have even more parameters. Therefore, a novel wrapper-based technique employing GEO is needed for feature selection since it uses fewer tuning parameters and is hence computationally less expensive. While meta-heuristic algorithms have been successful in choosing the most beneficial features that improve classification[25] performance, they often suffer from high exploration, which increases the likelihood that the optimization algorithm will miss the overall best solution and remain at local optima. Moreover, some current feature selection models have improved performance, but only for specific datasets, highlighting the need for hybrid approaches that achieve a balance amongst exploration and exploitation. However, despite the success of the existing models, the no free lunch (NFL) in optimization suggests that there is no single optimization method that can address all optimization problems.

Therefore, a new hybrid approach is necessary to address FS issues. The complexity of FS is a significant challenge, and many meta-heuristic approaches have been suggested to manage it. However, due to the increased exploitation of these methods, there is a greater chance that the optimization algorithm may miss the global optimum solution and become trapped at local optima. Furthermore, only a small number of the currently available feature selection strategies have demonstrated superior performance for all datasets, highlighting the need for hybrid feature selection approaches that balance local and global search. Despite the current models' accomplishments, the NFL in optimization revealed that no method can identify the ideal feature sets for all datasets in the field of feature selection. A chaotic map that merges with POA is apparently advantageous for feature selection since POA is a recently found robust optimization strategy and has demonstrated its usefulness in hybrid schemes. Therefore, a novel hybrid approach [2] that incorporates chaotic map with POA is needed to address the complexity of feature selection and achieve optimal results for different datasets.

### **1.3 CONTRIBUTIONS**

- Perform BGEO-TVFL algorithm by using KNN (K-Nearest Neighbour) model to find the no of features selected and accuracy on medical datasets
- Perform BWOA algorithm by using KNN (K-Nearest Neighbour) model to find the no of features selected and accuracy on medical datasets
- Perform BGWO algorithm by using KNN (K-Nearest Neighbour) model to find the no of features selected and accuracy on medical datasets
- Perform ACO algorithm by using KNN (K-Nearest Neighbour) model to find the no of features selected and accuracy on medical datasets
- Perform ABC algorithm by using KNN (K-Nearest Neighbour) model to find the no of features selected and accuracy on medical datasets
- At last we conclude that BGEWO-TVFL algorithm is best compared to remaining all algorithms.

### **1.4 Overview of Binary Golden Eagle Optimizer**

The Binary Golden Eagle Optimizer (BGEO)[14] is an advanced optimization algorithm inspired by the foraging and hunting behavior of golden eagles, designed specifically for solving binary optimization problems. Its efficiency stems from its ability to mimic the eagle's strategic hunting patterns and adaptability. In BGEO, solutions are represented as binary vectors, where each bit can take a value of 0 or 1. This binary representation makes the algorithm computationally efficient and suitable for discrete optimization problems, such as feature selection in machine learning and engineering design.

#### **Key Features and Mechanisms:**

##### **1. Exploration (Spiral Paths):**

During the exploration phase, BGEO mimics the eagle's spiral flight patterns, enabling a broad search across the solution space. This phase is crucial for identifying promising regions that may contain the optimal solution. The wandering behavior in the early stages reflects the eagle's natural tendency to explore widely, ensuring a high likelihood of discovering diverse solutions.

## **2. Exploitation (Linear Paths and Attacks):**

As the algorithm progresses, it simulates the eagle's focused straight-line attack behavior. This stage refines the search by concentrating on the most promising areas identified during exploration, increasing the precision and convergence rate toward the global optimum. The shift from wandering to attacking is a strategic transition, balancing exploration and exploitation.

## **3. Dynamic Mode Switching:**

A notable feature of BGEO is its flexibility in switching attack modes. This dynamic adaptability allows the algorithm to avoid local optima by adjusting its strategy based on the fitness landscape. Each "flight" or iteration involves reassessing the environment and dynamically altering the search mode, demonstrating its robustness in handling complex optimization problems.

## **4. Information Gathering:**

Golden eagles in nature rely not only on their instincts but also on observing and learning from other eagles' foraging activities. Similarly, BGEO incorporates a collaborative mechanism where information about high-quality solutions is shared among search agents.

This aspect enhances the algorithm's ability to converge efficiently by leveraging collective knowledge.

## **5. Fitness-Based Regeneration:**

BGEO employs a fitness function to evaluate and regenerate binary vectors. Solutions with higher fitness are more likely to influence the regeneration process, ensuring that the search progresses toward optimal solutions. This probability-based selection mechanism maintains a balance between global search and local refinement.

## 2. LITERATURE SURVEY

ShenkaiGu and others. [9] Proposed a new feature selection method which combined a Competitive Swarm Optimizer (CSO) and a clustering method to reduce the computation cost. This method applied for selected features of the concentration decreased, resulting in better classification performance. However, the authors noted that this approach still needs to be optimized to achieve the best balance between feature selection and classification performance.

R. K. Agarwal and others. [10] Proposed a wrapper feature selection method based on Quantum Wolf Algorithm Optimization (Quantum WOA), which combined quantum concept with WOA to enhance exploitation and detection efficiency Crossover-optimized mutation using quantum rotary gate and quantum bit representation uses. Q-bits improved with operators, and overall performance increased but this approach was limited to solving single- objective optimization problems and did not extend to continuous improvement or multi- objective optimization problems

Xian-fang Song et al. [11] present an efficient feature selection method that combines bare bone particle swarm optimization (BBPSO) and mutual information to solve high-dimensional feature selection problems Method was initiated by a swarm initialization model that used label correlation to enhance convergence. Local search operators two deletion and complement operators were developed to determine feature relevance and redundancy and improve exploitation performance Furthermore, an adaptive flip mutation operator was added to help particles escape local optima and find optimal solutions. The BBPSO method was able to obtain the feature subset with improved performance, although it took longer to complete the feature selection task compared to the other methods

M. S. Uzer, N. Yilmaz, and O. Inan [12] proposed an Artificial Bee Swarm-Based Feature Selection Method aimed at improving the diagnosis of liver diseases, cirrhosis, and diabetes using a Support Vector Machine (SVM) model. The method employed a swarm-based optimization technique to identify the most relevant features while excluding redundant and irrelevant ones, enhancing the model's predictive performance.

J. Too and S. Mirjalili (2021) [13] introduced an improved Binary Dragonfly Algorithm (BDA) to address the challenge of extracting local optima in wrapper-based feature selection optimization problems. The proposed algorithm, called Hyper Learning Binary Dragonfly

Algorithm (HLBDA), improves the search behavior by incorporating a hyper learning method to dynamically balance exploration and exploitation during optimization.

Got et al. (2021), [26] had suggested a new combined filter-wrapper FS scheme using the WOA. The suggested approach was a multi-objective algorithm in which fitness functions for a wrapper and a filter were simultaneously optimized. An exhaustive evaluation of this method against seven well-known algorithms on 12 benchmark database shows the effectiveness of our technique. The suggested approach examines the most recent Guided Population Archive WOA (GPAWOA) system and merges a filter and wrapper model into a single system in an effort to capitalize on the advantages of each model. As a result, during the training phase, the filter and wrapper assessment criteria were KNN classifier and mutual information. Additionally, the hyperbolic tangent function was used to give GPAWOA the ability to handle discrete issues

Abdel-Basset et al. (2020),[19] had suggested a new GWO method combine two-phase mutation address FS for classification issues based on wrapper approaches. Two-phase mutation improves the algorithm's ability to be exploited. Decrease sum of selected features while keeping good classification accuracy was the goal of the first mutation phase. The second mutation phase's goal was to try to boost the classification accuracy by adding more detailed traits. The two-phase mutation can be accomplished with a low probability because the mutation phase can be time-consuming.

Hussien et al. (2019), had proposed the best feature subset for the dimensionality reduction and classification process was chosen using a novel binary variant of WOA. An S-shaped sigmoid transfer function serves as the foundation of the new methodology. In order to convert a free location of the whale to their corresponding binary solutions, the feature selection problem must be addressed. This transformation was accomplished by applying an S-shaped transfer function to each dimension, which determines the chance of varying the elements of the position vectors from 0 to 1 and vice versa. This forces the search agents to travel in a binary space. Application of the K-NN classifier ensures that the chosen characteristics were the pertinent ones. The suggested BWOA-S was compared to the native one using a set of criteria across eleven various database. The outcomes demonstrated that the new method performs noticeably well in locating the ideal feature.

H. Rao et al. (2019) [18] proposed a feature selection method combining the Artificial Bee Colony (ABC) algorithm with the Gradient Boosting Decision Tree (GBDT) model to enhance classification accuracy. ABC optimizes feature selection by identifying relevant attributes, while GBDT improves predictive performance through iterative decision tree refinement. This hybrid approach outperforms traditional methods like PCA and RFE by eliminating strict data distribution assumptions and offers better stability than Genetic Algorithms (GA) and Particle Swarm Optimization (PSO). The study highlights the effectiveness of integrating swarm intelligence with machine learning-based boosting, making it useful for medical applications like breast cancer detection. Since your Bio-CAD project focuses on optimized feature selection, this research aligns with your use of Binary Golden Eagle Optimization (BGEO) and Time-Varying Flight Length (TVFL). Incorporating boosting techniques like GBDT could further improve the accuracy and efficiency of your model in medical diagnosis.

E. Emay et al. (2016)[19] introduced a feature selection method using Binary Grey Wolf Optimization (BGWO), inspired by the hierarchical leadership and hunting behavior of grey wolves. This approach efficiently selects the most relevant features while minimizing redundancy, leading to improved classification accuracy. Compared to traditional feature selection techniques like Principal Component Analysis (PCA) and Recursive Feature Elimination (RFE), BGWO offers better flexibility by avoiding strict data distribution assumptions. It also demonstrates superior performance over other bio-inspired optimization algorithms such as Genetic Algorithms (GA) and Particle Swarm Optimization (PSO) in terms of convergence speed and stability. The study highlights the advantages of swarm intelligence for feature selection, particularly in medical and pattern recognition applications. Implementing BGWO or similar bio-inspired techniques could further enhance classification accuracy and computational efficiency, making your model more effective for real-world medical diagnosis.

A. Mohammadi-Balani et al. (2021)[20] introduced the Golden Eagle Optimizer (GEO), a nature- inspired metaheuristic algorithm based on the hunting strategy of golden eagles. This optimization technique mimics the intelligent flight and target-tracking behavior of eagles to explore and exploit the search space efficiently. GEO has demonstrated strong performance in solving complex

optimization problems by balancing exploration and exploitation, making it competitive with other metaheuristic algorithms such as Genetic Algorithms (GA), Particle Swarm Optimization (PSO), and Grey Wolf Optimization (GWO). Unlike traditional optimization methods, GEO adapts dynamically to problem constraints, leading to improved convergence speed and accuracy. The study highlights its applications in feature selection, engineering optimization, and machine learning tasks. Since your Bio-CAD project utilizes Binary Golden Eagle Optimization (BGEO) and Time- Varying Flight Length (TVFL) for medical dataset classification, this research directly aligns with your approach. Leveraging GEO's adaptability and efficiency could further enhance feature selection and classification accuracy, making it a valuable component in optimizing medical diagnosis models.

M. Abdel-Basset et al. (2020)[21] proposed an enhanced feature selection method by fusing the Grey Wolf Optimizer (GWO) with a two-phase mutation strategy to improve exploration and exploitation. The hybrid approach refines the traditional GWO by incorporating a mutation mechanism that prevents premature convergence and enhances search efficiency. This method outperforms standard feature selection techniques like Principal Component Analysis (PCA) and Recursive Feature Elimination (RFE) by dynamically adjusting the balance between exploration and exploitation, leading to more optimal feature subsets. Compared to other metaheuristic algorithms such as Genetic Algorithms (GA) and Particle Swarm Optimization (PSO), the fusion approach offers better stability and adaptability, making it highly effective in classification tasks. The study highlights its practical applications in medical diagnosis and machine learning, where selecting the most relevant features is crucial for accurate predictions. , this research aligns with your optimization-driven feature selection approach. Integrating similar hybrid strategies, such as mutation-enhanced metaheuristics, could further improve the accuracy and efficiency of your model, making it more robust for real-world medical diagnosis.

Q. Al-Tashi et al. (2019)[22] introduced a hybrid feature selection method using Binary Grey Wolf Optimization (BGWO) combined with other optimization strategies to enhance performance. The proposed approach improves the traditional Grey Wolf Optimization (GWO) by refining the search process for selecting the most relevant features while reducing redundancy. This hybrid method effectively balances exploration and exploitation, leading to better convergence and classification accuracy. Compared to conventional techniques like Principal Component Analysis

(PCA) and Recursive Feature Elimination (RFE), it provides greater flexibility and adaptability. Additionally, it outperforms other bio-inspired algorithms such as Genetic Algorithms (GA) and Particle Swarm Optimization (PSO) in terms of efficiency and stability. The study demonstrates its effectiveness in medical and machine learning applications, where optimal feature selection is crucial for improving predictive models. Since your Bio-CAD project focuses on medical dataset classification using Binary Golden Eagle Optimization (BGEO) and Time-Varying Flight Length (TVFL), this research aligns with your approach. Implementing hybrid optimization strategies, such as BGWO, could further enhance the feature selection process, improving both accuracy and computational efficiency for medical diagnosis.

S. Mirjalili et al. (2020) introduced the Whale Optimization Algorithm (WOA), a nature-inspired metaheuristic optimization technique based on the bubble-net hunting strategy of humpback whales. The algorithm mimics the encircling, spiraling, and searching behaviors of whales to balance exploration and exploitation, making it highly effective for solving complex optimization problems. WOA has been successfully applied in various fields, including feature selection, engineering design, and machine learning. Compared to traditional methods like Principal Component Analysis (PCA) and Recursive Feature Elimination (RFE), WOA dynamically adapts to the problem space, leading to improved optimization performance. It also competes well with other bio-inspired algorithms such as Particle Swarm Optimization (PSO) and Grey Wolf Optimization (GWO) in terms of convergence speed and solution quality. The study highlights its application in designing photonic crystal filters, demonstrating its effectiveness in handling high-dimensional optimization problems. Since your Bio-CAD project focuses on medical dataset classification using Binary Golden Eagle Optimization (BGEO) and Time-Varying Flight Length (TVFL), WOA's adaptive search mechanism could be explored to further enhance feature selection and classification accuracy, potentially improving the efficiency of medical diagnostic models.

## 3. SYSTEM ANALYSIS

### 3.1 EXISTING SYSTEMS

Traditional feature selection methods, such as the Competitive Swarm Optimizer (CSO) with clustering, aim to reduce computational costs and improve classification accuracy. This method enhances feature selection by removing redundant data, leading to better model performance. However, it still struggles to find the optimal balance between feature selection and classification efficiency. Many existing algorithms face challenges such as high computational complexity, slow convergence, and suboptimal feature subset selection. These limitations affect overall system performance, making it necessary to develop a more adaptive and efficient feature selection approach to enhance medical data classification.

The Quantum Wolf Algorithm Optimization (Quantum WOA) is a wrapper-based feature selection method that integrates quantum computing concepts with WOA to enhance exploitation and detection efficiency. However, this approach is limited to single-objective optimization and does not support continuous improvement or multi-objective optimization problems. Due to these constraints, it struggles to adapt dynamically to different datasets and optimize multiple conflicting objectives, highlighting the need for more flexible and adaptive feature selection techniques in medical data classification.

The Bare Bone Particle Swarm Optimization (BBPSO) with Mutual Information is an efficient feature selection method designed to handle high-dimensional datasets. It improves convergence speed using a swarm initialization model that leverages label correlation. Additionally, local search operators help refine feature selection by reducing redundancy. To further enhance performance, an adaptive flip mutation operator allows particles to escape local optima and achieve better global solutions. However, despite its improved feature selection capability, BBPSO requires longer processing time, making it less efficient for real-time applications and large-scale medical datasets.

The Artificial Bee Swarm-Based Feature Selection Method was developed to enhance the diagnosis of liver diseases, cirrhosis, and diabetes using a Support Vector Machine (SVM) model.

This approach utilizes swarm-based optimization to select the most relevant features while eliminating redundant and irrelevant ones, thereby improving classification accuracy. By leveraging collective intelligence and adaptive searching, the method enhances the predictive performance of the SVM classifier. However, despite its effectiveness, swarm-based techniques often suffer from slow convergence and high computational cost, making them less efficient for large-scale medical datasets.

The improved Binary Dragonfly Algorithm (BDA) was introduced to tackle the challenge of extracting local optima in wrapper-based feature selection optimization problems. The proposed algorithm, known as Hyper Learning Binary Dragonfly Algorithm (HLBDA), enhances search behavior by incorporating a hyper learning method that dynamically balances exploration and exploitation during optimization. This approach helps the algorithm avoid premature convergence and improves feature selection efficiency. However, despite these improvements, the method may still face challenges in handling high-dimensional datasets and computational complexity, which can affect real-time performance in medical data classification.

The combined filter-wrapper feature selection scheme using the Whale Optimization Algorithm (WOA) was proposed as a multi-objective approach, optimizing both wrapper and filter fitness functions simultaneously. During the training phase, the filter and wrapper assessment criteria were based on the KNN classifier and mutual information. Additionally, the hyperbolic tangent function was incorporated to enable GPAWOA to handle discrete feature selection problems effectively. While this approach demonstrated strong performance across multiple benchmark datasets, its computational complexity and optimization overhead may limit its scalability for large medical datasets.

Various optimization-based feature selection methods have been proposed to enhance classification accuracy. Techniques such as the Binary Dragonfly Algorithm (BDA), Artificial Bee Swarm Optimization, and Grey Wolf Optimizer (GWO) address challenges related to local optima and feature selection. Some approaches integrate wrapper and filter methods, like the Guided Population Archive WOA (GPAWOA), to optimize feature subsets effectively. Others, such as the Hyper Learning Binary Dragonfly Algorithm (HLBDA), improve search behavior using adaptive strategies. While these methods enhance performance, they often involve high

computational complexity and may face difficulties in handling large-scale medical datasets, affecting their real-time efficiency.

Hussien et al. (2019) proposed a binary variant of the Whale Optimization Algorithm (WOA) for feature selection and dimensionality reduction. This approach utilizes an S-shaped sigmoid transfer function to map continuous whale positions into binary values, enabling effective feature selection. The transformation process determines the probability of switching elements between 0 and 1, ensuring that search agents operate in a binary space. The selected features were evaluated using the K-NN classifier to retain the most relevant ones. The proposed Binary Whale Optimization Algorithm with Sigmoid function (BWOA-S) was tested on eleven different datasets, and results showed significant improvement in identifying the optimal feature subset compared to the standard WOA.

H. Rao et al. (2019) proposed a feature selection method combining the Artificial Bee Colony (ABC) algorithm with the Gradient Boosting Decision Tree (GBDT) model to improve classification accuracy. ABC identifies relevant features, while GBDT enhances predictive performance through iterative learning. This approach outperforms traditional methods like PCA and RFE by avoiding strict data distribution assumptions and offers better stability than Genetic Algorithms (GA) and Particle Swarm Optimization (PSO). The integration of swarm intelligence with boosting techniques makes it effective for medical applications like breast cancer detection. This method aligns with feature selection techniques like Binary Golden Eagle Optimization (BGEO) and Time-Varying Flight Length (TVFL) used in your Bio-CAD project

E. Emry et al. (2016) proposed a feature selection method using Binary Grey Wolf Optimization (BGWO), inspired by grey wolves' hierarchical hunting behavior. This approach enhances classification accuracy by selecting relevant features while minimizing redundancy. It outperforms traditional methods like PCA and RFE by avoiding strict data distribution assumptions and surpasses GA and PSO in convergence speed and stability. The study highlights BGWO's effectiveness in medical applications like breast cancer detection. This aligns with your Bio-CAD project, where incorporating BGWO could further improve feature selection and classification accuracy.

## 3.2 Disadvantages of Existing System

- Developing optimization methods that efficiently handle high-dimensional data while maintaining accuracy.
- Designing optimization-based feature selection methods that are robust to imbalanced datasets.
- Ensuring that the selected features and the decision-making process are interpretable to domain experts.
- Developing Ulti-objective optimization algorithms that balance competing objectives.
- Reducing the computational cost without compromising performance.
- Creating optimization-based approaches that are robust [15] to noisy and redundant features.

## 3.3 PROPOSED SYSTEMS

We use BGEO-TVFL (Binary Golden Eagle Optimizer-Time Varying Flight Length) algorithm to solve Feature Selection problem. The Binary Golden Eagle Optimization (BGE) algorithm is then applied to generate a number of binary vectors, each representing a possible feature subset, using a fitness function based on a K-Nearest Neighbour (KNN) classifier time- varying flight length (2010). TVFL parameter inserted is used BGE-algorithm is to vary the length of each flight (iteration) based on the current population health, and enable the algorithm to navigate the required complex areas and search areas This procedure can have found more solutions by increasing the flight length as fitness improves, also By decreasing it as fitness stops, the algorithm is able to intensify the search around promising regions BGE algorithm perform runs for a specified number of generations or until the stopping criterion is met, after which the top features are selected based on the frequency of appearance from the final population , evaluated using different metrics, including accuracy, precision, recall, F1-score, and area under the receiver operating characteristic curve (AUC-ROC) The performance of the trained classifier is evaluated using both training and testing datasets, whose distribution is obtained comprehensive power analysis This new approach is expected to improve classification performance and robustness. The BGEO solution is described by the e dimensional vector Y where Y represents as a continuous value vector then the new locations of the attributes is updated by the following equation:

$$y^{u+1} = y^u + \Delta y^u \quad i$$

Where,  $\Delta y^u$  represents the  $j$  phase

$$\Delta y^u = \rightarrow. TF \xrightarrow{t=1} \xrightarrow{j} + \rightarrow. r. \xrightarrow{v} \xrightarrow{\frac{A}{2l}} \xrightarrow{j} \rightarrow. \xrightarrow{D} (2)$$

$j$  vector of features at iteration

selected items and maximum classification accuracy) [47], and the fitness function is defined as below:

$$Fit(Y) = \beta \lambda (Y) + \eta \frac{|Y|}{|M|} \quad (3)$$

Where,  $\lambda S (Y)$  denotes the classification error rate,  $|M|$  means that the transfer function is used to determine the optimal result. Furthermore, the transfer function is considered to have an integral and important function in the derivation of the BGEO\_TVFL algorithm changing the transfer function can change the performance and the total number of features in data set, defines the selected number of features by  $Y$ , denoting the weights  $\beta$  and  $\eta$  classification error rate and selection coefficient,  $\beta$

1. The relationship between  $\beta$  and  $\eta$  indicating  $\eta \propto 1/\beta$ . It is well preserved. Last but not least are the steps for the proposed BGEO-TVFL described here in the number of items originally produced. A set of data is obtained and represented by. Next, a 2D array of size  $M$ ,  $e$  is determined,

### **Algorithm 1: Time-varying flight length and binary Golden Eagle optimizer**

set the feature population to start.

Identify the fitness function.

Set up the position initially.

Set each memory's population to zero.

Set  $F$  to initial, update  $F$  for each iteration  $v$ , and for each GE  $j$ .

Choose a trait at random from the populations' memory

Calculate AV

If the length of AV is not zero

Establish your CV.

Utilizing equation (1) calculate the step vector  $Ay$ .

Equation (2) is used to update the new position. Equation (3) is used to convert to binary.

Calculate the fitness function at the new position by utilizing equation (4) Once the fitness function has outperformed the memory location in GE  $j$

Replace the current location with the one stored in GE  $j$ 's memory

End

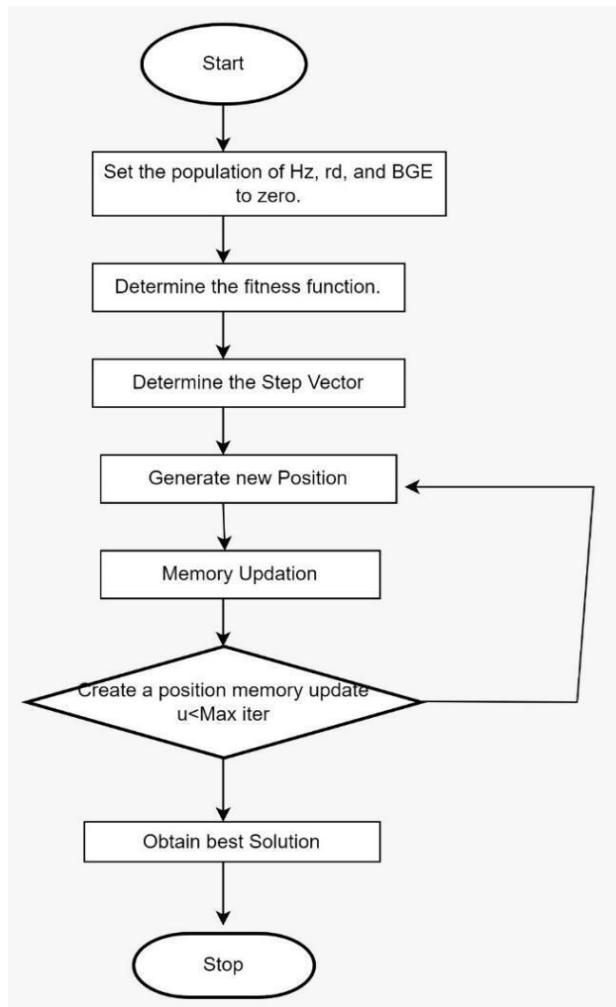
End

End

End

and its value is defined as 0 or 1. But M defines the number of attributes when e is equal to the number of objects found in the dataset. The feature vector is designed to look like any row. Which is based on the feature vector is presented for the purpose. Moreover, the memory fitness of each GE or feature can be calculated and initialized. Now, the process is just repeated for Max\_iter multiple times. For all iterations, the state of the GE was updated, then the validated position, which indicates eligibility for the relevant position, was updated accordingly and finally the rest of the GE was updated accordingly. The BGEO-TVFL pseudo code is given as below

The ability of the BGEO-TVFL algorithm to efficiently explore the solution space and adapt to the changing fitness states is due to the unique combination of plane dynamics and binary encoding - It can be highly adjustable, so as not to hit the local optimal, and converge to global optimum Binary encoding scheme makes the algorithm more representable and efficient in the solution space, in order to search the search space more efficiently and also reduces the risk of early convergence f, the power of the algorithm is less from experience and adapting to other cases makes it a promising optimization technique for real-world applications such as complex system optimization, machine learning, and industrial design. Time varying flight length and binary Golden Eagle (GE) optimizer is a new algorithm that combines the concept of time varying flight length and Golden Eagle optimizer to solve complex optimization problems Feature that the algorithm starts with a fitness function -Puts the population of the location is then initialized, and the number of people in each memory is set to zero. The algorithm repeats each iteration, v, randomly selecting an element from the remaining j population for each Golden Eagle (GE). The algorithm calculates the Average Velocity (AV) and if it is not zero, updates the location using Equations (2) and (3), binary transforms the new location with Equation (3) and then evaluates the fitness function at the new location , if and If the memory space of GE j is exceeded, the current space is replaced by the space reserved in the memory of GE j This process is repeated for each iteration and GE, allowing the algorithm to adapt to it corresponds to the optimization problem and converges to the optimal solution



**Fig. 2. Flowchart for BGEO-TVFL**

The BGEO-TVFL algorithm initializes the Hz, rd, and BGE values to zero, then determines the fitness function, which evaluates the goodness of each solution in the population. Calculating the Step Vector, which determines the direction and magnitude of each search term inside Generates new locations, searching the solution space efficiently. The Memory Update mechanism is used to update the position memory, which stores the best solution found so far, and repeat this process for the maximum number of iterations, a Mixite Updates Specified, iterations continue until the maximum number of iterations is reached. Finally, the algorithm finds and prunes the best solution found in the optimization process, using binary coding to find the solution reach the most challenging areas through the principles of aviation energy and optimize the changing state of fitness.

## **3.4 FEASIBILITY STUDY**

A feasibility study is a structured analysis used to assess the practicality and viability of a proposed project or system. It evaluates various factors such as technical, operational, economic, legal, and schedule feasibility to determine whether the project can be successfully implemented.

### **3.4.1 Technical Feasibility**

- **Algorithm Efficiency:** The BGEO-TVFL model outperforms existing optimization algorithms in feature selection and classification accuracy.
- **Computational Requirements:** The algorithm requires moderate computational resources, as indicated by varying computational times for different datasets
- **Scalability:** The model has been tested on multiple datasets with different feature sizes and sample counts, demonstrating adaptability to various medical datasets.

### **3.4.2 Operational Feasibility**

- **Ease of Use:** The system integrates well with medical datasets and provides interpretable results for healthcare professionals.
- **Automation and Decision Support:** The feature selection process reduces complexity, aiding doctors in diagnosis without needing extensive manual feature engineering.
- **Reliability:** Consistent classification accuracy across datasets suggests robust performance in real-world medical applications.

### **3.4.3 Economic Feasibility**

- Cost-Benefit Analysis:
- Reduces computational cost by selecting fewer features while maintaining high accuracy.
- Potential savings in medical diagnostics by improving classification accuracy, reducing unnecessary tests.
- Infrastructure Requirements: The model runs efficiently on general computational setups without excessive hardware investment.

### **3.4.4 Economic Feasibility**

**Patient Data Privacy:** The model must comply with data protection laws (e.g., HIPAA, GDPR) when handling medical datasets.

- **Medical Standards Compliance:** AI-based models must meet regulatory standards for clinical decision support systems.

### **3.4.5 Schedule Feasibility**

- **Project Timeline:**

- Model Development & Testing: 3 months
- Evaluation & Optimization: 2 months
- Deployment & Integration: 3 months

- **Potential Challenges:**

- Potential Challenges: Larger datasets may require extended preprocessing.

## **3.5 COCOMO MODEL**

It involves implementing and evaluating multiple machine learning models for feature selection and classification. COCOMO can be used to estimate the effort, cost, and time required for software development based on factors like project complexity, team experience, and system requirements.

### **COCOMO Model Overview**

COCOMO estimates project effort (in person-months) using the formula:

$$E = a(KLOC)^b$$

where:

E = Effort (in person-months) KLOC = Thousands of lines of code a, b = Model-specific

constants The estimated time (T) required is:

$$T = c(E)^d \text{ where:}$$

T = Time in months

c, d = Model-specific constants

COCOMO has three modes based on project complexity:

1. Organic – Small, simple projects with experienced teams
  2. Semi-Detached – Moderate complexity with mixed experience levels.
  3. Embedded – Complex projects with strict constraints
- Estimating Effort for Your Project

For your BGEO-TVFL and KNN-based classification system, we can estimate development effort as follows:

1. Determine KLOC: Measure the total lines of code in your feature selection and classification implementation.
2. Choose the COCOMO Model Type:

If it is a research prototype with limited complexity, Organic mode fits.

If integrating multiple algorithms and handling large medical datasets, Semi-Detached is better.

3. Estimate Effort and Cost:

Using the Basic COCOMO model, for Semi-Detached mode:  $E = 3.0 \times (KLOC)^{1.12} \times T = 2.5 \times (E)^{0.35}$

]

If your project is around 10 KLOC, then:

$$E = 3.0 \times (10)^{1.12} = 33.1 \text{ person-months}$$

$T = 2.5 \times (33.1)^{0.35} = 8.1$  months

]

## 4. SYSTEM REQUIREMENTS

### 4.1 Software Requirements:

- i. **OperatingSystem** : Windows11,64-bitOperatingSystem
- ii. **Coding Language** : Python
- iii. **Python distribution** :Google Colab Pro, Flask
- iv. **Browser** :Any Latest Browser like Chrome

To develop and execute the BGEO-TVFL and KNN-based Medical Data Classification System, various software components are required, including an operating system, programming languages, libraries, frameworks, and development tools. The system is compatible with Windows 10/11 (64- bit), Linux (Ubuntu 20.04, CentOS), and macOS (Big Sur or later), ensuring a stable platform for running machine learning models and optimization algorithms.

The primary programming language for implementation is Python (3.8 or later), as it provides robust libraries for machine learning, optimization, and data processing. If a web-based user interface (UI) is required, Angular can be used for frontend development, while Flask or Django can be utilized for backend integration to process medical datasets. Several key libraries are essential for system functionality, including NumPy, Pandas, and Scikit-learn for data manipulation and classification, SciPy for scientific computing, and Matplotlib/Seaborn for result visualization. For feature selection and optimization, DEAP or PyCharm can be used to implement BGEO-TVFL. Performance evaluation metrics such as accuracy, precision, recall, F1-score, and AUC-ROC can be computed using sklearn.metrics. If deep learning is incorporated in future enhancements, TensorFlow or PyTorch can be considered. In terms of storage, databases such as SQLite, MySQL, or PostgreSQL can be used to store datasets, results, and reports for easy retrieval and comparison.

The development process can be streamlined using Jupyter Notebook, PyCharm, or VS Code, while Google Colab can be an alternative for cloud-based execution with GPU support. Git and GitHub are recommended for version control and collaboration.

If deployment is needed, cloud platforms like AWS EC2, Google Cloud, or Microsoft Azure can host the system, while Flask/Django APIs can facilitate interaction between the machine learning model and a web-based UI. Bootstrap or Tailwind CSS can be used for UI styling in case a web frontend is developed.

## 4.2 Hardware Requirements:

- i. **System Type** :64-bit operating system, x64-based processor
- ii. **Cache memory** :4MB(Megabyte)
- iii. **RAM** :16GB(gigabyte)
- iv. **Hard Disk** :8GB

To efficiently execute the BGEO-TVFL and KNN-Based Medical Data Classification System, a suitable hardware setup is essential to handle medical datasets, feature selection, and classification tasks. The system requires a 64-bit processor with a minimum clock speed of 2.5 GHz, preferably an Intel Core i5/i7 or AMD Ryzen 5/7 for smooth computation. Since feature selection and classification involve high computational loads, having at least 8GB of RAM is recommended, while 16GB or more is preferable for large medical datasets like the Thyroid dataset (3153 samples, 23 features).

A dedicated GPU (such as NVIDIA GeForce GTX 1650 or RTX 2060 and above) can significantly enhance processing speed, especially if deep learning models or parallel computing techniques are incorporated in the future. The system should have a minimum of 500GB SSD storage, with 1TB SSD or HDD recommended to store multiple datasets, classification results, and logs efficiently. A high-resolution monitor (1080p or higher) is beneficial for analyzing data visualizations, classification reports, and performance metrics.

For cloud-based execution or remote deployment, the system can leverage AWS, Google Cloud, or Azure Virtual Machines, which provide scalable computing power based on demand. If a web-based UI is developed, a dedicated server with a high-speed internet connection is required to ensure smooth real-time processing of medical data. The hardware setup should also support external storage devices for backup and secure handling of sensitive medical information.

### **4.3 SOFTWARE DESCRIPTION:**

The BGEO-TVFL and KNN-Based Medical Data Classification System is designed to optimize feature selection and improve classification accuracy for medical datasets. The system is implemented using Python (3.8 or later) due to its extensive libraries for machine learning, optimization, and data processing. If a web-based interface is required, Angular is used for frontend development, while Flask or Django serves as the backend framework for processing medical datasets.

For data handling and preprocessing, NumPy and Pandas are used to manipulate datasets, clean missing values, normalize data, and encode categorical values. The Scikit-learn library is integrated for applying the K-Nearest Neighbour (KNN) classifier, which predicts medical conditions based on selected features. The SciPy library is leveraged for advanced mathematical computations and optimization processes. The core feature selection method, Binary Golden Eagle Optimization with Time-Varying Flight Length (BGEO-TVFL), is implemented using optimization libraries like DEAP or PySwarm, ensuring efficient selection of relevant features.

To evaluate system performance, sklearn.metrics provides classification metrics such as accuracy, precision, recall, F1-score, and AUC-ROC. Matplotlib and Seaborn are utilized for visualizing classification results, feature selection efficiency, and performance comparisons with other optimization algorithms such as BWOA, BGWO, ACO, and ABC.

If database storage is required, the system supports SQLite, MySQL, or PostgreSQL to store medical datasets, classification results, and historical records for future reference. Development and execution are carried out using Jupyter Notebook, PyCharm, or VS Code, with Google Colab available for cloud-based execution with GPU acceleration. Git and GitHub ensure version control and collaboration.

For deployment, AWS, Google Cloud, or Microsoft Azure provide scalable cloud computing services, enabling real-time processing and remote access. If the system includes a web-based UI, Bootstrap or Tailwind CSS can be used for responsive styling, enhancing user.

## 5. SYSTEM DESIGN

### 5.1 SYSTEM ARCHITECTURE:

#### 5.1.1 Input Datasets

The system utilizes five medical datasets with different feature sizes and classification challenges:

- **Breast Cancer Dataset:** Predicts whether a tumor is malignant or benign (binary classification, 9 features).
- **Heart Disease Dataset:** Diagnoses cardiovascular diseases (multi-class classification, 13 features).
- **Diabetes Dataset:** Identifies whether a patient has diabetes (binary classification, 8 features).
- **Iris Dataset:** Classifies iris flowers into three species (multi-class classification, 4 features).
- **Thyroid Dataset:** Diagnoses thyroid conditions (binary classification, 23 features).

Each dataset contains multiple instances with varying numbers of features that must be optimized for better classification accuracy.

#### 5.1.2 Data Preprocessing

Data preprocessing is a crucial step to clean and prepare raw data for feature selection and classification. The following steps are performed:

- **Handling Missing Values:** Imputation techniques such as mean/mode substitution or removal of incomplete records are applied.
- **Feature Scaling:** Normalization or standardization methods like Min-Max scaling are used to bring all features to a similar scale.
- **Data Encoding:** Categorical data (if any) is converted into numerical format using label encoding or one-hot encoding.
- **Data Splitting:** The dataset is divided into training (80%) and testing (20%) subsets to evaluate model performance.

### **5.1.3 Feature Selection (BGEO-TVFL Algorithm)**

The **Binary Golden Eagle Optimization (BGEO)** algorithm is implemented to select the most relevant features while reducing redundancy. It works in the following manner:

- **Initialization:** A population of binaryvectors (representing feature subsets) is generated.
- **Flight Length Variation (TVFL):** The flight length dynamically adjusts during iterations based on the search space to balance exploration and exploitation.
- **Fitness Evaluation:** A fitness function evaluates the qualityof selected features using a KNN- based classification error metric.
- **Selection and Updating:** Features with higher fitness are retained, and unimportant features are discarded.

BGEO-TVFL improves feature selection by dynamically refining the search space, leading to higher classification accuracy with fewer selected features.

### **5.1.4 Classification Model (KNN Algorithm)**

The **K-Nearest Neighbour (KNN)** algorithm [16] is used for classification, leveraging the optimal feature set obtained from BGEO-TVFL. The steps include:

- **Distance Calculation:** Computes the Euclidean distance between test samples and training samples.
- **Nearest Neighbor Selection:** Identifies the k-nearest data points based on distance.
- **Majority Voting:** Assigns the class label based on the majorityof k-neighbors.
- **Performance Metrics:** The classifier's effectiveness is evaluated using accuracy, precision, recall, F1-score, and AUC-ROC.

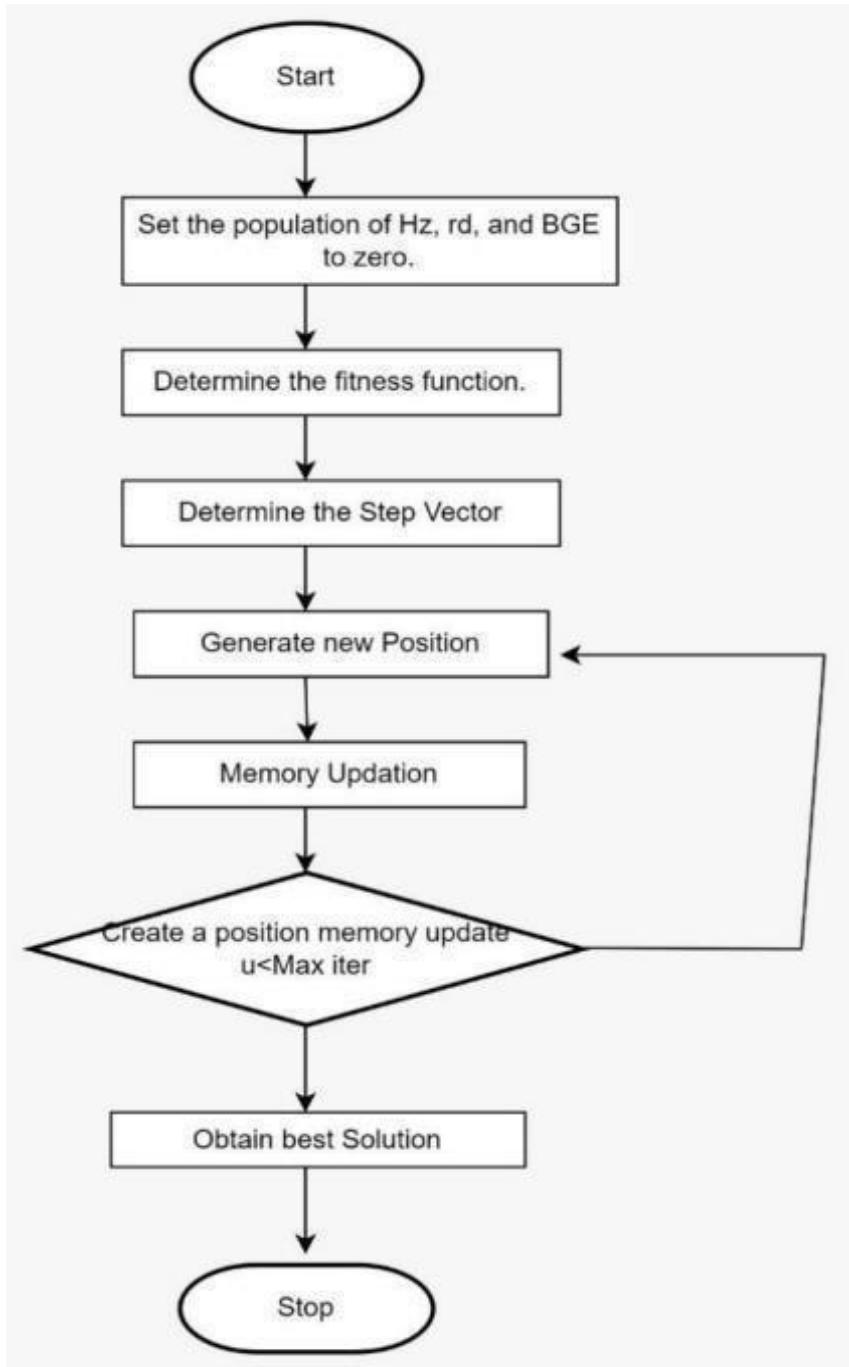
### **5.1.5 Model Evaluation and Comparison**

The proposed system is compared against other optimization algorithms, including BWOA (Binary Whale Optimization Algorithm), BGWO (Binary Grey Wolf Optimization), ACO (Ant Colony Optimization), and ABC (Artificial Bee Colony Optimization).

Key evaluation metrics include:

- **Classification Accuracy:** Measures how well the selected features contribute to correct predictions.

- **Computational Time:** Evaluates the time taken to process datasets.
- **Number of Selected Features:** Ensures minimal but relevant feature selection for better efficiency.



**Flowchart for BGEO-TVFL**

This flowchart represents an optimization algorithm workflow. It begins by initializing the population variables to zero. The fitness function is then determined to evaluate potential solutions. A step vector is calculated to guide movement within the solution space. The process continues until a maximum iteration limit is reached or a stopping criterion is met, at which point the best solution is obtained and the process stops.

space. New positions are generated iteratively, and memory is updated to store progress. The process continues until the maximum iteration condition is met. Finally, the best solution is obtained, and the algorithm terminates.

## 5.2 MODULES

### 5.2.1 Data Collection Module

**Objective:** Gather relevant medical datasets for feature selection and classification.

**Description:**

- **Dataset Selection:** Choose diverse medical datasets such as:
  - **Breast Cancer Dataset:** Used for detecting malignant and benign tumors.
  - **Heart Disease Dataset:** Determines the presence of cardiovascular diseases.
  - **Diabetes Dataset:** Identifies diabetic patients based on medical attributes.
  - **Iris Dataset:** Classifies iris flower species based on physical characteristics.
  - **Thyroid Dataset:** Detects thyroid conditions.
  - **Data Acquisition:** Datasets are collected from open-source repositories like Kaggle, UCI Machine Learning Repository, and medical research papers.
  - **Data Storage Format:** The datasets are stored in CSV, Excel, or relational databases to facilitate easy access and processing.

### 5.2.2 Data Preprocessing Module

**Objective:** Clean and standardize raw data to improve model performance.

**Description:**

- **Handling Missing Values:** Missing values are replaced using:
  - **Mean/Median Imputation:** For numerical attributes.
  - **Mode Imputation:** For categorical attributes.
- **Row Deletion:** If the missing data percentage is too high.
- **Feature Scaling:** Normalization techniques ensure all features have a consistent scale.
- **Min-Max Scaling:** Transforms data to a range (0,1).
- **Z-score Standardization:** Converts data into standard normal distribution.
- **Encoding Categorical Variables:** Converts non-numerical data into numeric format using:
  - **Label Encoding:** Assigns a unique integer to each category.
  - **One-Hot Encoding:** Creates separate binary columns for categorical values.

- **Data Splitting:** The dataset is divided into:
- **Training Set (80%):** Used to train the model.
- **Testing Set (20%):** Used to evaluate performance.

### **5.2.3 Feature Selection Module (BGEO-TVFL Algorithm)**

**Objective:** Select the most relevant features for improved classification accuracy.

**Description:**

- **Binary Golden Eagle Optimization (BGEO) Process:**
  - Generates a population of binary feature selection vectors.
  - Uses a fitness function based on classification error rate.
- **Time-Varying Flight Length (TVFL) Mechanism:**
  - Adapts the search space dynamically, preventing premature convergence.
  - Increases search length when the solution is suboptimal.
  - Reduces search length when nearing an optimal solution.
- **Feature Subset Selection:**
  - Irrelevant or redundant features are eliminated.
  - The model retains only the most informative attributes.
- **Optimization Objective:**
  - **Maximize Accuracy:** Select the best subset of features.
  - **Minimize Feature Count:** Reduce complexity while maintaining performance.

### **5.2.4 Classification Module (KNN Model)**

**Objective:** Classify medical conditions using the optimal feature set.

**Description:**

- **K-Nearest Neighbour (KNN) Classification Process:**
  - **ComputeDistance:** Calculates the Euclidean distance between test and training samples.
  - **Find k-Nearest Neighbors:** Identifies the closest data points based on distance.
  - **Majority Voting:** Determines the class label by selecting the most common class among neighbors.
- **Performance Metrics Evaluated:**
  - **Accuracy:** Measures correct classifications.
  - **Precision:** Evaluates the proportion of positive predictions that are correct.
  - **Recall (Sensitivity):** Identifies how well the model detects actual positive cases.

- **F1-score:** Balances precision and recall.

### **5.2.5 Model Evaluation and Comparison Module**

**Objective:** Compare the effectiveness of BGEO-TVFL with other optimization algorithms.

**Description:**

- **Baseline Algorithms for Comparison:**
- **Binary Whale Optimization Algorithm (BWOA).**
- **Binary Grey Wolf Optimization (BGWO).**
- **Ant Colony Optimization (ACO).**
- **Artificial Bee Colony Optimization (ABC).**
- **Comparison Parameters:**
- **Classification Accuracy:** Measures how well the model classifies medical data.
- **Number of Selected Features:** Determines efficiency in reducing dataset dimensionality.
- **Computational Time:** Evaluates how quickly each algorithm processes the data.
- **Visualization of Results:**
- Accuracy graphs comparing different models.
- Feature selection comparison tables.
- Radar charts representing performance across datasets.

### **5.2.6 System Integration Module**

**Objective:** Integrate all components into a functional system for end-to-end execution.

**Description:**

- **Backend Implementation:**
- Stores preprocessed datasets, selected features, and classification results.
- Uses Python with Flask or Django for API development.
- **Frontend Development (Optional):**
  - **Angular-based Web Interface** to visualize feature selection and classification results.
  - Displays accuracy, feature importance, and prediction results.
- **End-to-End Execution:**
  - Automates the entire process from **data input → preprocessing → feature selection → classification → evaluation.**

### **5.2.7 Deployment and Testing Module**

**Objective:** Deploy the system for real-world medical applications and validate model performance.

**Description:**

- **Testing on Multiple Datasets:** Ensures generalizability across medical datasets.
- **Hyperparameter Tuning:** Adjusts KNN parameters (k-value, distance metric) for better accuracy.
- **Deployment Options:**
- **Cloud Deployment:** Using AWS, Google Cloud, or Azure for real-time medical analysis.
- **Local Server Deployment:** Running the model on a hospital or research institution's internal system.

### **5.2.8 Report Generation Module**

**Objective:** Generate automated reports summarizing classification results.

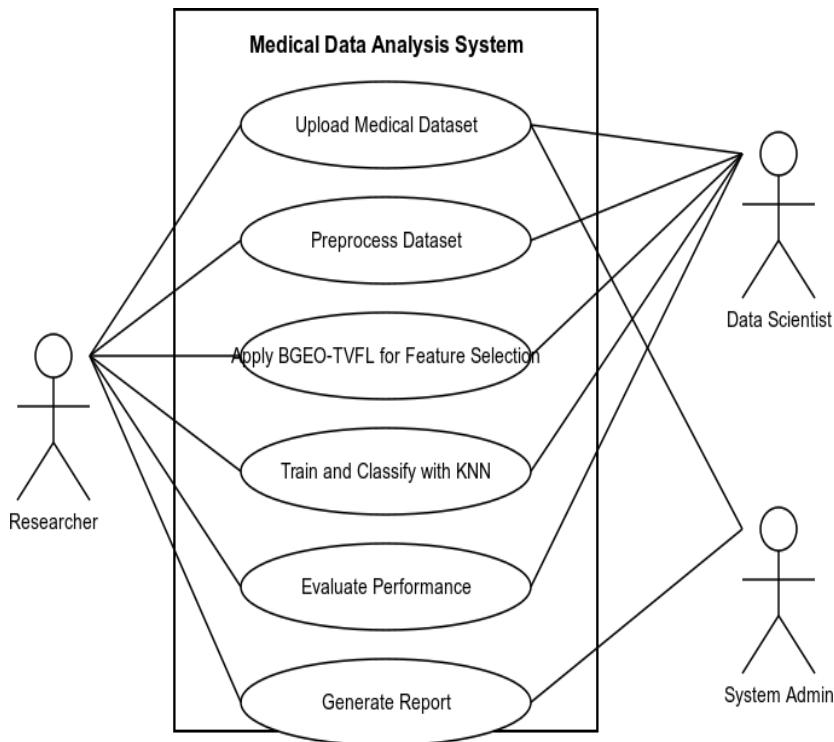
**Description:**

- **Automated Report Generation:**
  - Summarizes key findings of feature selection and classification.
  - Includes accuracy scores, selected features, and dataset insights.
- **Exporting Reports:**
  - **CSV Format:** For further data processing.
  - **PDF Format:** For easy sharing with healthcare professionals.
- **Graphical Representations:**
  - Accuracy comparison charts.
  - Feature selection visualization.
  - Confusion matrix for classification evaluation.

## 5.3 UML DIAGRAMS

### 5.3.1 USECASE DIAGRAM

A **Use Case Diagram** is a visual representation of how a system interacts with its users and other systems. It shows **who** (actors) can perform **what actions** (use cases) within a system.

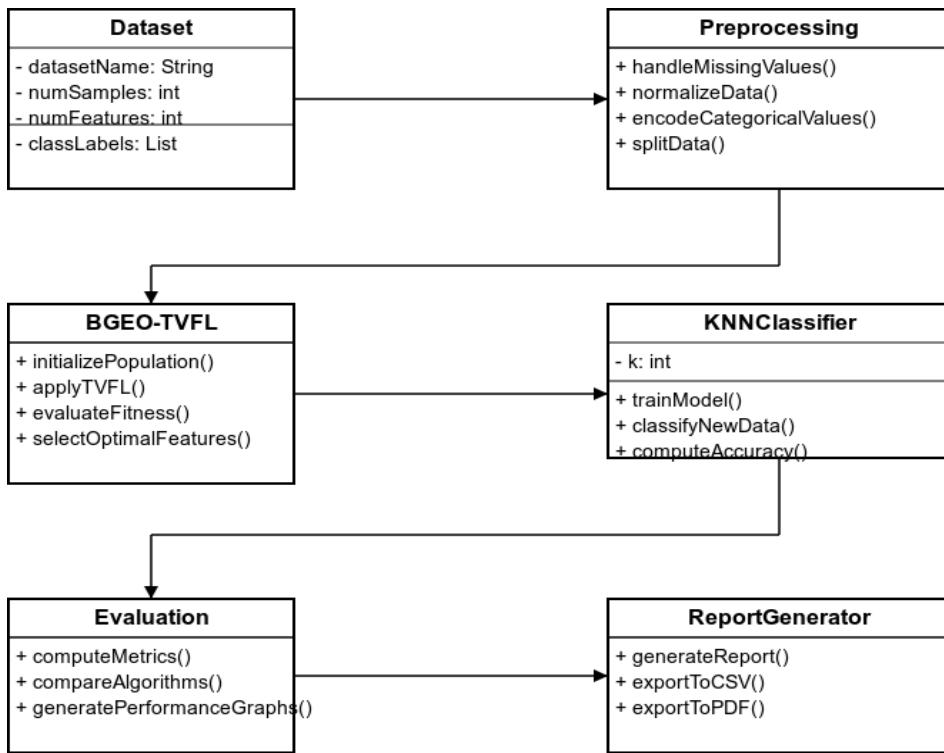


**Fig: 3Usecase Diagram**

### 5.3.2 CLASS DIAGRAM

A **Class Diagram** in UML is a visual representation of a system's **structure**. It shows:

- **Classes** (blueprints of objects)
- **Attributes** (data stored in a class)
- **Methods/Functions** (actions a class can perform)
- **Relationships** (how classes interact with each other).



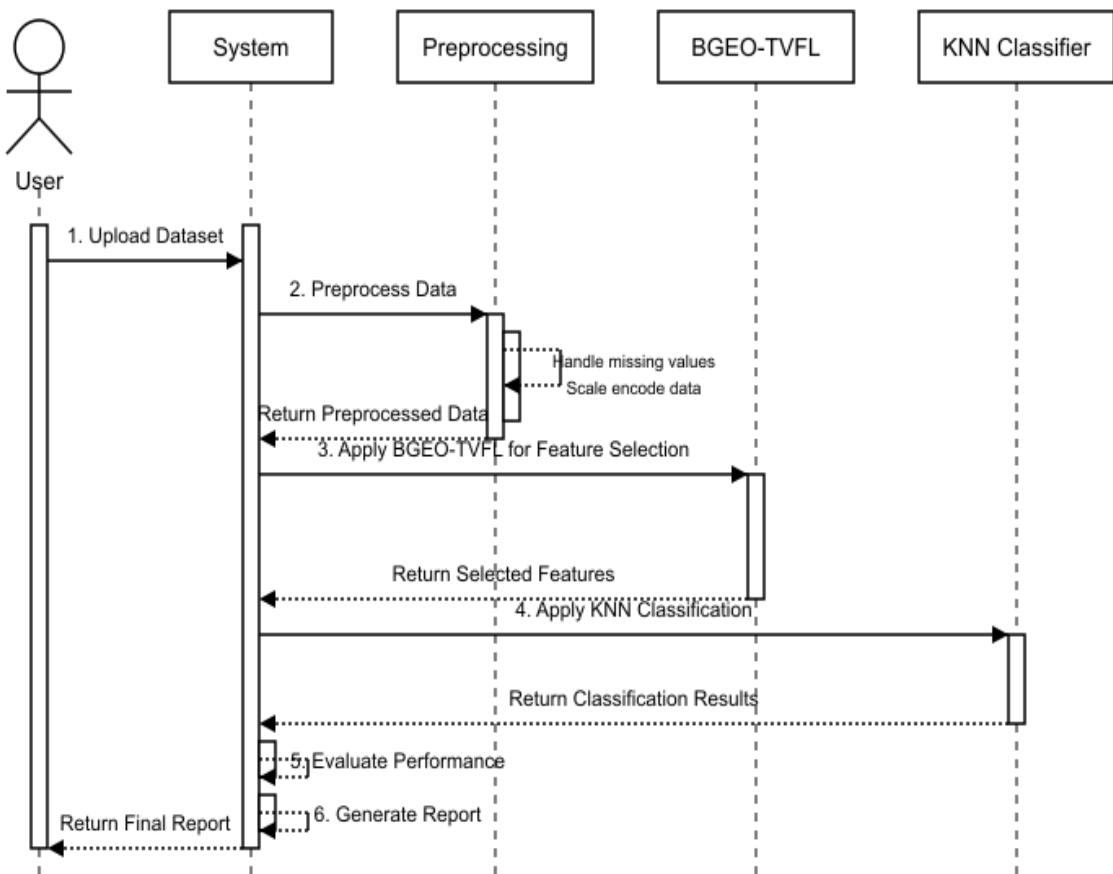
**Fig:4 Class Diagram**

### 5.3.3 SEQUENCE DIAGRAM

A **Sequence Diagram** in UML is a visual representation of how objects in a system interact over time. It shows:

- **Objects** (participants in the system)
- **Lifelines** (the duration an object exists in a scenario)
- **Messages** (interactions between objects, shown as arrows)
- **Activation Bars** (when an object is actively processing a request)

### Medical Data Analysis System - Sequence Diagram



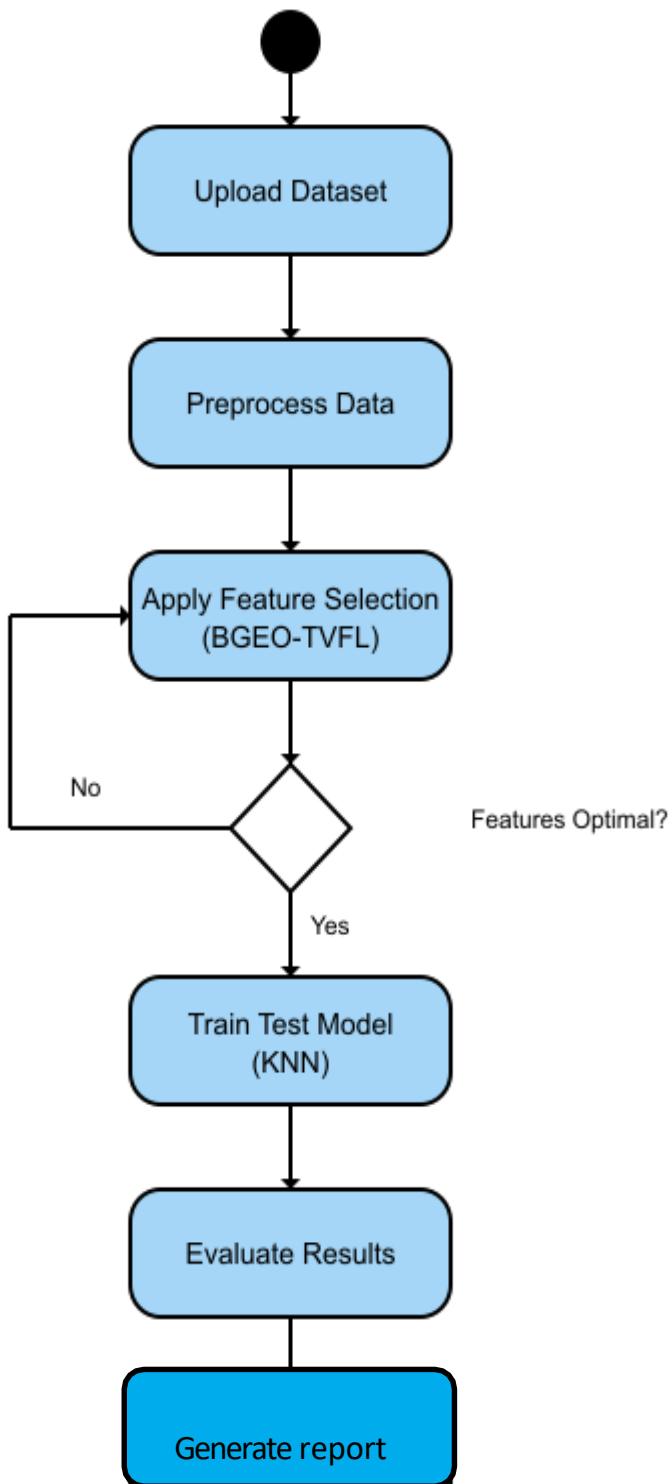
**Fig:5 Sequence Diagram**

#### **5.3.4 ACTIVITY DIAGRAM**

An **Activity Diagram** in UML is a visual representation of the **workflow or process flow** in a system. It shows:

- **Activities** (tasks or steps in the process)
- **Decision Points** (conditional branches)
- **Start & End Nodes** (beginning and completion of the process)
- **Arrows** (flow of control from one activity to another)

## Medical Data Analysis System - Activity Diagram



**Fig:6 Activity Diagram**

## 6. IMPLEMENTATION

### 6.1 Model Implementation

Feature selection for BGEO-TVFL

```
#KNN for breast_cancer.csv using BGEO-TVFL
import numpy as np
Import pandas as pd
import time
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score

# Load the breast cancer dataset
data_path = "/content/drive/My Drive/AG4"

# Read the breast cancer data from the CSV file
breast_cancer_data = pd.read_csv(os.path.join(data_path, "breast_cancer.csv"))

# Assuming the last column is the target variable and the rest are features
X = breast_cancer_data.iloc[:, :-1].values # Use breast_cancer_data instead of data
y = breast_cancer_data.iloc[:, -1].values # Use breast_cancer_data instead of data

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Standardize the features if necessary
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# BGEO TVFL Parameters
num_eagles = 20
max_iter = 100
dimension = X_train.shape[1]

# Initialize eagles' positions
eagles_position = np.random.randint(2, size=(num_eagles, dimension))

# Fitness function: accuracy of K-Nearest Neighbors
def fitness_knn(selected_features):
if np.sum(selected_features) == 0: return 0
X_train_selected = X_train[:, selected_features == 1]
```

```

model=KNeighborsClassifier(n_neighbors=5)
model.fit(X_train_selected, y_train)
y_pred=model.predict(X_test_selected)
return accuracy_score(y_test, y_pred)

# BGEOTVFLOptimization
Loop best_solution = None
best_score = -1
worst_score =
1 all_scores =
[]
all_solutions=[]
average_solution = np.zeros(dimension)

# Measure start time
start_time=time.time()

for iteration in range(max_iter):
    iteration_best_score = 0
    iteration_best_solution = None

    for i in range(num_eagles):
        # Evaluate fitness of the current position
        score = fitness_knn(eagles_position[i])
        all_scores.append(score)
        all_solutions.append(eagles_position[i].copy())

    # Update best and worst score if
    score > best_score:
        best_score = score
        best_solution = eagles_position[i].copy()
    score < worst_score:
        worst_score = score

    # Update iteration best score and solution if
    score > iteration_best_score:
        iteration_best_score = score
        iteration_best_solution =
eagles_position[i].copy()

    # Apply the BGEOTVFL algorithm logic
    l=0.05 +(0.95 - 0.05) * (iteration / max_iter) if
    np.random.rand() < 0.5:
        new_position = best_solution - l * np.abs(best_solution - eagles_position[i]) else:
        new_position = best_solution + l * np.abs(best_solution - eagles_position[i])

    # Convert position to binary based on thresholding
    new_position = 1 / (1 + np.exp(-new_position)) new_position =
np.where(new_position >= 0.5, 1, 0)

```

```

# Update the position
eagles_position[i]=new_position

# Calculate the average solution for this iteration
average_solution += np.array(iteration_best_solution)

# Print iteration results
print(f"Iteration {iteration + 1}: Best Score (Accuracy) = {iteration_best_score:.4f}, Best Solution = {iteration_best_solution}")
print(f"Overall Best Score (Accuracy) = {best_score:.4f}, Overall Best Solution = {best_solution}")

# Measure end time
end_time=time.time()
total_time = end_time - start_time

# Calculate average fitness
average_fitness=np.mean(all_scores)

# Calculate average solution
average_solution /= max_iter

# Print final results
print("Final Best Solution:", best_solution)
print("Final Best Accuracy:", best_score)
print("Worst Accuracy:", worst_score)
print("Average Accuracy:", average_fitness)
print("Average Best Solution:", average_solution)
print("Total Time Taken:", total_time, "seconds")

#KNN forbreast_cancer.csv using BWOA
import numpy as np
import pandas as pd
import time
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score

# Load the breast cancer dataset
data_path = "/content/drive/My Drive/AG4"

# Read the breast cancer data from the CSV file
breast_cancer_data = pd.read_csv(os.path.join(data_path, "breast_cancer.csv"))

# Assuming the last column is the target variable and the rest are features
X=breast_cancer_data.iloc[:, :-1].values # Use breast_cancer_data instead of
data y = breast_cancer_data.iloc[:, -1].values # Use breast_cancer_data

```

instead of data

```
# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

#Standardizethe
features
scaler =
StandardScaler()
X_train=scaler.fit_transform(X_train)

X_test = scaler.transform(X_test)

#BWOAParameters
num_whales = 20
max_iter = 100
dimension = X_train.shape[1]

# Initialize whales' positions
whales_position = np.random.randint(2, size=(num_whales, dimension))

# Fitness function: accuracy of K-Nearest
Neighboursdeffitness_knn(selected_features):
if np.sum(selected_features) ==
0: return 0
X_train_selected = X_train[:, selected_features == 1]
X_test_selected = X_test[:, selected_features == 1]

model=KNeighborsClassifier(n_neighbors=5)
model.fit(X_train_selected, y_train)
y_pred=model.predict(X_test_selected)
return accuracy_score(y_test, y_pred)

#BWOAOptimization
Loop best_solution =
None best_score = -1
worst_score = 1
all_scores = []
average_solution = np.zeros(dimension)

#Measurestarttime
start_time =
time.time()

for iteration in
range(max_iter):
iteration_best_score = 0
```

```

for i in range(num_whales):
    # Evaluate fitness of the current position
    score = fitness_knn(whales_position[i])
    all_scores.append(score)

    # Update best and worst score
    if score > best_score:
        best_score = score
        best_solution=whales_position[i].copy()
    if score < worst_score:
        worst_score = score

    # Update iteration best score and solution
    if score > iteration_best_score:
        iteration_best_score = score
        iteration_best_solution =
        whales_position[i].copy()

    # Apply the BWOA algorithm logic
    l=-1+(1-(-1))*(iteration / max_iter)
    p = np.random.rand()
    if p < 0.5:
        if np.abs(l) < 0.5:
            D=np.abs(2 * np.random.rand() * best_solution - whales_position[i])
            new_position = best_solution - l * D
        else:
            random_whale=whales_position[np.random.randint(0, num_whales)] D=
            np.abs(2 * np.random.rand() * random_whale - whales_position[i])
            new_position = random_whale - l * D
        else:
            distance_to_leader = np.abs(best_solution - whales_position[i])
            new_position = distance_to_leader * np.exp(l) * np.cos(2 * np.pi * l) + best_solution

    # Convert position to binary based on thresholding
    new_position = 1 / (1 + np.exp(-new_position))
    new_position=np.where(new_position >= 0.5, 1, 0)

    # Update the position
    whales_position[i]=new_position

    # Calculate the average solution for this iteration
    average_solution +=
    np.array(iteration_best_solution)

# Print iteration results
print(f"Iteration {iteration + 1}: Best Score (Accuracy) = {iteration_best_score:.4f}, Best
Solution = {iteration_best_solution}")
print(f"Overall Best Score (Accuracy) = {best_score:.4f}, Overall Best Solution =

```

```

{best_solution}")

#Measureendtime
end_time =
time.time()
total_time = end_time - start_time

#Calculateaveragefitness
average_fitness =
np.mean(all_scores)

# Calculate average solution
average_solution /= max_iter

# Print final results
print("Final Best Solution:", best_solution)

print("Final Best Accuracy:", best_score)
print("Worst Accuracy:", worst_score)
print("Average Accuracy:", average_fitness)
print("Average Best Solution:", average_solution)
print("Total Time Taken:", total_time, "seconds")

#KNN forbreast_cancer.csv using BGWO
import numpy as np
import pandas as pd
import time
from sklearn.model_selection import
train_test_split from sklearn.preprocessing
import StandardScaler from sklearn.neighbors
import KNeighborsClassifier from
sklearn.metrics import accuracy_score

# Load the breast cancer dataset
data_path = "/content/drive/My Drive/AG4"

# Read the breast cancer data from the CSV file
breast_cancer_data = pd.read_csv(os.path.join(data_path, "breast_cancer.csv"))

# Assumingthe last column is the target variable and the rest are
features X = breast_cancer_data.iloc[:, :-1].values
y = breast_cancer_data.iloc[:, -1].values

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

#Standardizethe
features scaler =

```

```

for i in range(num_wolves):
    StandardScaler()
    X_train =
        scaler.fit_transform(X_train)
    X_test=scaler.transform(X_test)

#BGWOParameters
num_wolves = 20
max_iter = 100
dimension = X_train.shape[1]

# Initialize wolves' positions
wolves_position = np.random.randint(2, size=(num_wolves, dimension))

# Fitness function: accuracy of K-Nearest
Neighborsdeffitness_knn(selected_features):
if np.sum(selected_features) == 0:
    return 0
X_train_selected=X_train[:,selected_features == 1]
X_test_selected = X_test[:, selected_features == 1]

model = KNeighborsClassifier()
model.fit(X_train_selected, y_train)
y_pred =
model.predict(X_test_selected)return
accuracy_score(y_test, y_pred)

# BGWO Optimization Loop
alpha_position =
np.zeros(dimension)alpha_score=
-1
beta_position =
np.zeros(dimension) beta_score=-
1
delta_position =
np.zeros(dimension)delta_score=
-1
all_scores = []
average_solution = np.zeros(dimension)

#Measurestarttime
start_time =
time.time()

for iteration in range(max_iter):
iteration_best_score = 0
iteration_best_solution = None

```

```

# Evaluate fitness of the current position
score = fitness_knn(wolves_position[i])
all_scores.append(score)

# Update alpha, beta, and delta wolves
if score > alpha_score:
    delta_score = beta_score
    delta_position = beta_position.copy()
    beta_score = alpha_score
    beta_position = alpha_position.copy()
    alpha_score = score
    alpha_position = wolves_position[i].copy()
elif score > beta_score:
    delta_score = beta_score
    delta_position = beta_position.copy()
    beta_score = score
    beta_position = wolves_position[i].copy()
    alpha_score = alpha_position.copy()
    alpha_position = wolves_position[i].copy()

# Update iteration best score and solution
if score > iteration_best_score:
    iteration_best_score = score
    iteration_best_solution = wolves_position[i].copy()

# Update positions of
# wolves for i in
# range(num_wolves):
for j in range(dimension):
    r1 = np.random.rand() r2
    =np.random.rand() A1 =
    2 * r1 - 1
    C1 = 2 * r2
    D_alpha = np.abs(C1 * alpha_position[j] - wolves_position[i][j]) X1
    =alpha_position[j] - A1 * D_alpha

    r1 = np.random.rand() r2
    =np.random.rand() A2 =
    2 * r1 - 1
    C2 = 2 * r2
    D_beta = np.abs(C2 * beta_position[j] - wolves_position[i][j]) X2
    =beta_position[j] - A2 * D_beta
    r1 = np.random.rand() r2
    =np.random.rand() A3 =
    2 * r1 - 1
    C3 = 2 * r2

```

```

D_delta=np.abs(C3 * delta_position[j] - wolves_position[i][j]) X3
= delta_position[j] - A3 * D_delta

wolves_position[i][j] =(X1 + X2 + X3)/ 3

# Applysigmoid function to update position wolves_position[i][j]
= 1 / (1 + np.exp(-wolves_position[i][j]))

# Convert to binary
wolves_position[i][j] = np.where(wolves_position[i][j] >= 0.5, 1, 0)

# Calculate the average solution for this iteration
if iteration_best_solution is not None: # Check if a solution was found in this iteration
average_solution += np.array(iteration_best_solution)

# Print iteration results
print(f"Iteration {iteration + 1}: Best Score (Accuracy) = {iteration_best_score:.4f}, Best
Solution = {iteration_best_solution}")
print(f"Overall Best Score (Accuracy) = {alpha_score:.4f}, Overall Best Solution =
{alpha_position}")

#Measureendtime
end_time =
time.time()
total_time = end_time - start_time

#Calculateaveragefitness
average_fitness =
np.mean(all_scores)

# Calculate average solution
average_solution /= max_iter

# Print final results
print("Final Best Solution:", alpha_position)
print("Final Best Accuracy:", alpha_score)
print("Worst Accuracy:", min(all_scores))
print("Average Accuracy:", average_fitness)
print("Average Best Solution:",
average_solution) print("Total Time Taken:",
total_time, "seconds")

#KNN for breast_cancer.csv using
ACO import numpy as np

```

```

import pandas as pd
pd import time
from sklearn.model_selection import train_test_split from sklearn.preprocessing import StandardScaler from sklearn.neighbors import KNeighborsClassifier from sklearn.metrics import accuracy_score

# Load the breast cancer dataset
data_path = "/content/drive/MyDrive/AG4"
beta = 1

# Initialize pheromone levels
pheromone =
np.ones(dimension)

# Fitness function: accuracy of k-Nearest Neighbours
def fitness_knn(selected_features):
if np.sum(selected_features) == 0:
return 0
X_train_selected = X_train[:, selected_features == 1]
X_test_selected = X_test[:, selected_features == 1]

model = KNeighborsClassifier(n_neighbors=5)
model.fit(X_train_selected, y_train)
y_pred = model.predict(X_test_selected)
return accuracy_score(y_test, y_pred)

# ACO Optimization Loop
best_solution =
np.zeros(dimension)
best_score
=-1
all_scores = []
average_solution = np.zeros(dimension)

# Measure start time
start_time =
time.time()

for iteration in range(max_iter):
iteration_best_score = 0
iteration_best_solution = None
solutions = []

for ant in range(num_ants):
solution = np.zeros(dimension)
for i in range(dimension):
probability = (pheromone[i] * alpha) * ((1.0 / (i + 1)) * beta) if

```

```

np.random.rand() < probability:
    solution[i] = 1

    score = fitness_knn(solution)
    all_scores.append(score)
    solutions.append((solution,score))

    if score > iteration_best_score:
        iteration_best_score = score
        iteration_best_solution = solution.copy()
    if score > best_score:
        best_score = score
        best_solution = solution.copy()

#Update pheromone
levels for i in
range(dimension):
    pheromone[i] *=(1 - evaporation_rate)
for solution, score in solutions:
    if solution[i] == 1:
        pheromone[i] += score

#Calculate the average solution for this iteration
average_solution +=
np.array(iteration_best_solution)

#Print iteration results
    print(f"Iteration {iteration + 1}: Best Score (Accuracy) = {iteration_best_score:.4f}, Best
Solution = {iteration_best_solution}")
    print(f"Overall Best Score (Accuracy) = {best_score:.4f}, Overall Best Solution =
{best_solution}")

#Measure end time
end_time =
time.time()
total_time = end_time - start_time

#Calculate average fitness
average_fitness =
np.mean(all_scores)

#Calculate average solution
average_solution /= max_iter

#Print final results
print("Final Best Solution:", best_solution)
print("Final Best Accuracy:", best_score)
print("Worst Accuracy:", min(all_scores))

```

```

print("Average Accuracy:", average_fitness)
print("Average Best Solution:",
average_solution) print("Total Time Taken:",
total_time, "seconds")
#KNN forbreast_cancer.csv using
ABC import numpy as np
import pandas as pd
from sklearn.model_selection import
train_test_split from sklearn.preprocessing
import StandardScaler from sklearn.neighbors
import KNeighborsClassifier from
sklearn.metrics import accuracy_score

# Load the breast cancer dataset
data_path = "/content/drive/My Drive/AG4"

# Read the breast cancer data from the CSV file
breast_cancer_data = pd.read_csv(os.path.join(data_path, "breast_cancer.csv"))

# Assumingthe last column is the target variable and the rest are
features X = breast_cancer_data.iloc[:, :-1].values
y=breast_cancer_data.iloc[:, -1].values

# Split data into training and test sets
X_train,X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Standardize the features
X_train =
scaler.fit_transform(X_train)
scaler=StandardScaler()
X_test = scaler.transform(X_test)

#ABCParameters
num_beans = 20
max_iter = 100
num_features = X_train.shape[1]

# Initialize population of bees with random feature selections
population = np.random.randint(2, size=(num_beans,
num_features))

#Fitness function: KNN accuracy
def
fitness_function(selected_features)
: if np.sum(selected_features) == 0:
return 0
X_train_selected=X_train[:, selected_features == 1]
X_test_selected = X_test[:, selected_features == 1]

```

```

model = KNeighborsClassifier()
model.fit(X_train_selected, y_train)
y_pred =
model.predict(X_test_selected) return
accuracy_score(y_test, y_pred)
# ABC Optimization
loop best_solution =
None best_accuracy
= -1
avg_best_solution = np.zeros(num_features)
best_fitness = -1
worst_fitness=
float('inf')
avg_fitness = 0

for iteration in range(max_iter):
    # Evaluate fitness for each bee's position
    fitness_scores = np.array([fitness_function(population[i]) for i in range(num_bees)])

    # Find best solution and update metrics
    iteration_best_index = np.argmax(fitness_scores)
    iteration_best_solution = population[iteration_best_index]
    iteration_best_fitness = fitness_scores[iteration_best_index]

    if iteration_best_fitness > best_fitness:
        best_solution = iteration_best_solution
        best_accuracy = iteration_best_fitness
        best_fitness = iteration_best_fitness

    avg_best_solution += iteration_best_solution
    avg_fitness += iteration_best_fitness

    if iteration_best_fitness < worst_fitness:
        worst_fitness = iteration_best_fitness

    # ABC: Update population based on
    # fitness for i in range(num_bees):
    # Generate new solution for each bee
    candidate_solution = population[i]
    neighbor_bee_index = np.random.choice(np.delete(np.arange(num_bees), i))
    neighbor_solution = population[neighbor_bee_index]

    if fitness_scores[i] < fitness_scores[neighbor_bee_index]:
        # Perform onlooker bee phase
        feature_to_flip = np.random.choice(np.where(candidate_solution != neighbor_solution)[0])
        candidate_solution[feature_to_flip] = 1 - candidate_solution[feature_to_flip]

```

```

if fitness_scores[i] < best_fitness:
#Perform scout bee phase
num_flips = np.random.randint(1, num_features)
random_features = np.random.choice(num_features, num_flips, replace=False)
candidate_solution[random_features] = 1 - candidate_solution[random_features]

population[i]=candidate_solution

# Print iteration results
print(f"Iteration {iteration + 1}: Best Accuracy= {iteration_best_fitness:.4f}")

#Calculate average best solution and fitness
avg_best_solution /= max_iter
avg_fitness /= max_iter

# Print final results
print("\nFinal Best Solution:", best_solution)
print("Final Best Accuracy:", best_accuracy)
print("Average Best Solution:",
avg_best_solution) print("Best Fitness:", best_fitness)
print("Worst Fitness:", worst_fitness)

print("Average Fitness:", avg_fitness)

```

## 6.2 CODING

```
# prompt: display all files in the directory AG4

import os

from google.colab import drive

import pandas as pd

drive.mount('/content/drive')

# Change the path to point to your AG4 dataset folder on Google Drive.

data_path = "/content/drive/My Drive/AG4"

# List all files in the AG4 directory

files = os.listdir(data_path)

#print the files

print(files)
```

### Preprocessing

Read and Print `breas_cancer.csv`

```
# prompt: read and print breast_cancer.csv

import pandas as pd

# Change the path to point to your AG4 dataset folder on Google Drive.

data_path = "/content/drive/My Drive/AG4"

# Read the breast cancer data from the CSV file

breast_cancer_data = pd.read_csv(os.path.join(data_path, "breast_cancer.csv"))

# Print the breast cancer data

print(breast_cancer_data)

Read the null values in breast_cancer.csv

# prompt: read the null values in breast_cancer.csv

# Read the breast cancer data from the CSV file
```

```

breast_cancer_data=pd.read_csv(os.path.join(data_path,"breast_cancer.csv"))

# Count the number of null values in each column

null_values=breast_cancer_data.isnull().sum()

# Print the number of null values

print(null_values)

#prompt: read the null values in breast_cancer.csv

Perform Statistical Analysis for breast_cancer.csv

# prompt: do statistical analysis for breast_cancer.csv

#Calculate descriptive statistics for the breast cancer data

descriptive_stats = breast_cancer_data.describe()

#Print the descriptive statistics

print(descriptive_stats)

#Calculate the correlation matrix for the breast cancer data

correlation_matrix = breast_cancer_data.corr()

#Print the correlation matrix

print(correlation_matrix)

Check the Outliers for breast_cancer.csv

#prompt: check the outliers for breast_cancer.csv

import seaborn as sns

import matplotlib.pyplot as plt

#Create a boxplot for each column in the breast cancer data

for column in breast_cancer_data.columns:

    sns.boxplot(x=breast_cancer_data[column])

    plt.show()

Remove the Outliers from breast_cancer.csv

# prompt: remove the outliers for breast_cancer.csv

```

```

import pandas as pd
from scipy import stats

# Read the breast cancer data from the CSV file
breast_cancer_data = pd.read_csv(os.path.join(data_path, "breast_cancer.csv"))

# Define a function to remove outliers using z-score
def remove_outliers(df, column, threshold=3):
    mean = df[column].mean()
    std = df[column].std()
    z_scores = (df[column] - mean) / std
    return df[(z_scores < threshold) & (z_scores > -threshold)] # Remove outliers for each column in the breast cancer data for
    column in breast_cancer_data.columns:
        breast_cancer_data = remove_outliers(breast_cancer_data, column)
# Print the breast cancer data after removing outliers
print(breast_cancer_data)

Print the Outliers after removing the outliers from breast_cancer.csv
# prompt: print the outliers graph after removing outliers from breast_cancer.csv

# Read the breast cancer data from the CSV file
breast_cancer_data = pd.read_csv(os.path.join(data_path, "breast_cancer.csv"))

# Define a function to remove outliers using z-score
def remove_outliers(df, column, threshold=3):
    mean = df[column].mean()
    std = df[column].std()
    z_scores = (df[column] - mean) / std
    return df[(z_scores < threshold) & (z_scores > -threshold)] # Remove outliers for each column in the breast cancer data
for column in breast_cancer_data.columns:

```

```

breast_cancer_data = remove_outliers(breast_cancer_data, column)

# Print the breast cancer data after removing outliers

print(breast_cancer_data)

# Create a boxplot for each column in the breast cancer data

for column in breast_cancer_data.columns:

    sns.boxplot(x=breast_cancer_data[column])

    plt.show()

Perform Correlation for breast_cancer.csv

# prompt: perform correlation for breast_cancer.csv

# Calculate the correlation matrix for the breast cancer data

correlation_matrix = breast_cancer_data.corr()

# Print the correlation matrix

print(correlation_matrix)

Perform separate Independent features and target variables for breast_cancer.csv

#prompt: do separate independent features and target variables for breast_cancer.csv

# Print the available columns to identify the correct target variable name

print(breast_cancer_data.columns)

# Separate independent features and target variable for breast cancer data

# Replace "Class" with the actual name if it differs

X = breast_cancer_data.drop("Glucose", axis=1)

y = breast_cancer_data["Glucose"] # Update with the actual target variable name

# Print the independent features

print(X)

# Print the target variable

print(y)

Perform normalization and standardization for breast_cancer.csv

```

```

#prompt: perform normalization and standardization for breast_cancer.csv

# Normalize the breast cancer data

from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()

X_normalized = scaler.fit_transform(X)

# Print the normalized data

print(X_normalized)

# Standardize the breast cancer data

from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()

X_standardized = scaler.fit_transform(X)

# Print the standardized data

print(X_standardized)

KNN for breast_cancer.csv using BGEO-TVFL

#KNN for breast_cancer.csv using BGEO-TVFL

import numpy as np

import pandas as pd

import time

from sklearn.model_selection import train_test_split

from sklearn.preprocessing import StandardScaler

from sklearn.neighbors import KNeighborsClassifier

from sklearn.metrics import accuracy_score

# Load the breast cancer dataset

data_path = "/content/drive/My Drive/AG4"

# Read the breast cancer data from the CSV file

breast_cancer_data = pd.read_csv(os.path.join(data_path, "breast_cancer.csv"))

```

```

# Assuming the last column is the target variable and the rest are features

X = breast_cancer_data.iloc[:, :-1].values # Use breast_cancer_data instead of data

y = breast_cancer_data.iloc[:, -1].values # Use breast_cancer_data instead of data #

Split the dataset into training and testing sets

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Standardize the features if necessary

scaler = StandardScaler()

X_train = scaler.fit_transform(X_train)

X_test = scaler.transform(X_test)

#BGEOTVFLParameters

num_eagles = 20

max_iter = 100

dimension = X_train.shape[1]

# Initialize eagles' positions

eagles_position = np.random.randint(2, size=(num_eagles, dimension))

# Fitness function: accuracy of K-Nearest Neighbors

def fitness_knn(selected_features):

    if np.sum(selected_features) == 0:

        return 0

    X_train_selected = X_train[:, selected_features == 1]

    X_test_selected = X_test[:, selected_features == 1]

    model = KNeighborsClassifier(n_neighbors=5)

    model.fit(X_train_selected, y_train)

    y_pred = model.predict(X_test_selected)

    return accuracy_score(y_test, y_pred)

# BGEOTVFL Optimization Loop

```

```

best_solution=None

best_score = -1

worst_score = 1

all_scores = []

all_solutions=[]

average_solution =np.zeros(dimension)

# Measure start time

start_time = time.time()

for iteration in range(max_iter):

    iteration_best_score = 0

    iteration_best_solution=None

    for i in range(num_eagles):

        # Evaluate fitness of the current position

        score = fitness_knn(eagles_position[i])

        all_scores.append(score)

        all_solutions.append(eagles_position[i].copy())

    #Update best and worst score

    if score > best_score:

        best_score = score

        best_solution=eagles_position[i].copy()

    if score < worst_score:

        worst_score = score

    # Update iteration best score and solution

    if score > iteration_best_score:

        iteration_best_score = score

        iteration_best_solution=eagles_position[i].copy()

```

```

# Apply the BGEOTVFL algorithm logic

l=0.05 + (0.95 - 0.05) * (iteration / max_iter)

if np.random.rand() < 0.5:

    new_position = best_solution - l * np.abs(best_solution - eagles_position[i])

else:

    new_position = best_solution + l * np.abs(best_solution - eagles_position[i])

# Convert position to binary based on thresholding

new_position = 1 / (1 + np.exp(-new_position))

new_position = np.where(new_position >= 0.5, 1, 0)

# Update the position

eagles_position[i] = new_position

# Calculate the average solution for this iteration average_solution +=

np.array(iteration_best_solution) # Print iteration results

print(f"Iteration {iteration + 1}: Best Score (Accuracy) = {iteration_best_score:.4f}, Best
      Solution =
{iteration_best_solution}")

print(f"Overall Best Score (Accuracy) = {best_score:.4f}, Overall Best Solution =
{best_solution}") # Measure end time

end_time = time.time()

total_time = end_time - start_time

# Calculate average fitness

average_fitness = np.mean(all_scores)

# Calculate average solution

average_solution /= max_iter

# Print final results

print("Final Best Solution:", best_solution)

```

```

print("Final Best Accuracy:", best_score)

print("Worst Accuracy:", worst_score)

print("Average Accuracy:", average_fitness)

print("Average Best Solution:", average_solution)

print("Total Time Taken:", total_time, "seconds")

KNN for diabetes.csv using BWOA

#KNN for diabetes.csv using BWOA

import numpy as np

import pandas as pd

import time

from sklearn.model_selection import train_test_split

from sklearn.preprocessing import StandardScaler

from sklearn.neighbors import KNeighborsClassifier

from sklearn.metrics import accuracy_score

# Load the breast cancer dataset

data_path = "/content/drive/My Drive/AG4"

# Read the breast cancer data from the CSV file

diabetes_data = pd.read_csv(os.path.join(data_path, "diabetes.csv"))

# Assuming the last column is the target variable and the rest are features

X = diabetes_data.iloc[:, :-1].values

y = diabetes_data.iloc[:, -1].values

# Split the dataset into training and testing sets

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Standardize the features

scaler = StandardScaler()

X_train = scaler.fit_transform(X_train)

```

```

X_test=scaler.transform(X_test)

# BWOA Parameters

num_whales = 20

max_iter = 100

dimension=X_train.shape[1]

# Initialize whales' positions

whales_position =np.random.randint(2, size=(num_whales, dimension))

# Fitness function: accuracy of K-Nearest Neighbours

def fitness_knn(selected_features):

    if np.sum(selected_features)==0:

        return 0

    X_train_selected=X_train[:, selected_features == 1]

    X_test_selected=X_test[:, selected_features == 1]

    model = KNeighborsClassifier(n_neighbors=5)

    model.fit(X_train_selected, y_train)

    y_pred=model.predict(X_test_selected)

    return accuracy_score(y_test, y_pred)

#BWOAOptimization Loop

best_solution = None

best_score = -1

worst_score=1

all_scores = []

average_solution=np.zeros(dimension)

# Measure start time

start_time = time.time()

for iteration in range(max_iter):

```

```

iteration_best_score = 0
iteration_best_solution = None
for i in range(num_whales):
    # Evaluate fitness of the current position
    score = fitness_knn(whales_position[i])
    all_scores.append(score)
    # Update best and worst score
    if score > best_score:
        best_score = score
        best_solution = whales_position[i].copy()
    if score < worst_score:
        worst_score = score
    # Update iteration best score and solution
    if score > iteration_best_score:
        iteration_best_score = score
        iteration_best_solution = whales_position[i].copy()
    # Apply the BWOA algorithm logic
    l = -1 + (1 - (-1)) * (iteration / max_iter)
    p = np.random.rand()
    if p < 0.5:
        if np.abs(l) < 0.5:
            D = np.abs(2 * np.random.rand() * best_solution - whales_position[i])
            new_position = best_solution - l * D
        else:
            random_whale = whales_position[np.random.randint(0, num_whales)]
            D = np.abs(2 * np.random.rand() * random_whale - whales_position[i])

```

```

new_position = random_whale - l * D

else:

    distance_to_leader = np.abs(best_solution - whales_position[i])

    new_position = distance_to_leader * np.exp(l) * np.cos(2 * np.pi * l) + best_solution

    # Convert position to binary based on thresholding

    new_position = 1 / (1 + np.exp(-new_position))

    new_position = np.where(new_position >= 0.5, 1, 0)

    # Update the position

    whales_position[i] = new_position

# Calculate the average solution for this iteration

average_solution += np.array(iteration_best_solution)

# Print iteration results

print(f"Iteration {iteration + 1}: Best Score (Accuracy) = {iteration_best_score:.4f}, Best Solution = {iteration_best_solution}")

print(f"Overall Best Score (Accuracy) = {best_score:.4f}, Overall Best Solution = {best_solution}")

# Measure end time

end_time = time.time()

total_time = end_time - start_time

# Calculate average fitness

average_fitness = np.mean(all_scores)

# Calculate average solution

average_solution /= max_iter

# Print final results

print("Final Best Solution:", best_solution)

print("Final Best Accuracy:", best_score)

print("Worst Accuracy:", worst_score)

```

```

print("Average Accuracy:", average_fitness)

print("Average Best Solution:", average_solution)

print("Total Time Taken:", total_time, "seconds")

```

## Frontend Code:

### Home.Html

```

<!DOCTYPE html>
<html lang="en">
<head>
    <meta charset="UTF-8">
    <meta name="viewport" content="width=device-width, initial-scale=1.0">
    <title>Medical Dataset Analysis</title>
    <link rel="stylesheet" href="{{ url_for('static', filename='style.css') }}>
    <linkrel="stylesheet" href="https://cdnjs.cloudflare.com/ajax/libs/font-
awesome/5.15.4/css/all.min.css">
    <style>
        .team {
            background: linear-gradient(to right, #2A9D8F, #E9C46A);
            text-align: center;
            padding: 40px 20px;
        }

        .team h2 {
            font-size: 2rem;
            margin-bottom: 20px;
            color: #333;
        }

        .team-list {
            list-style: none;
            padding: 0;
            margin: 0;
        }

        .team-list li {
            font-size: 1.2rem;
            margin-bottom: 10px;
            color: #555;
        }

    </style>
</head>
<body>
    <!-- PrimaryNavbar -->

```

```

<header class="navbar">
  <div class="container">
    <div class="logo-title">
      
      <h1>Classification and Feature Selection Method for Medical Datasets by BGEO-TVFL and KNN</h1>
    </div>
  </div>
</header>

<!-- Secondary Navigation Bar -->
<nav class="secondary-navbar">
  <div class="container">
    <a href="/about" class="nav-link">About</a>
    <a href="/flowchart" class="nav-link">Project Flowchart</a>
    <a href="/modelevaluationmetrics" class="nav-link">Model Evaluation Metrics</a>
  </div>
</nav>

<!-- Hero Section -->
<section class="hero">
  <div class="container hero-content">
    <div class="hero-text">
      <h2>Better Solutions for Your Medical Dataset Analysis</h2>
      <p>We are a team of experts leveraging advanced algorithms to optimize feature selection and classification.</p>
      <div class="hero-buttons">
        <a href="/upload" class="cta-btn">Get Started</a>
      </div>
    </div>
  </div>
</section>

<!-- Team Members Section -->
<section class="team">
  <div class="container">
    <h2>Meet Our Team</h2>
    <ul class="team-list">
      <li>Lakshmi Prasanna Bommu</li>
      <li>Swetha Lakshmi Paruchuri</li>
      <li>Sahithi Mellacheruvu</li>
    </ul>
  </div>
</section>

<!-- Footer -->
<footer class="footer">

```

```

<div class="container">
    <p>&copy; 2024 Medical Dataset Analysis Project. All rights reserved.</p>
</div>
</footer>
</body>
</html>

```

## App.py

```

from flask import Flask, render_template, request
import random # For generating random accuracies

(simulated) app = Flask(__name__)

# Simulating results for different methods
def get_accuracy_results(selected_method,
dataset_name): # Simulated accuracies for different
methods
methods = {
    "BGEO-TVFL": random.uniform(90, 95),
    "BWOA": random.uniform(85, 90),
    "BGWO": random.uniform(80, 85),
    "ACO": random.uniform(75, 80),
    "ABC": random.uniform(70, 75),
}

# Simulated time taken for execution
time_taken = random.uniform(1.0, 5.0) # Time in seconds

# Check if the selected method is
valid if selected_method in methods:
accuracy = methods[selected_method]
    return {selected_method: round(accuracy, 2)}, round(time_taken,
2) else:
# Return an error message if the method is invalid
return { "Error": f"Method '{selected_method}' not found"}, round(time_taken, 2)

@app.route(
 '/')
 def
home():
return render_template('home.html')

@app.route('/abo
ut')
def about():
return render_template('about.html')

```

```

@app.route('/upload',methods=['GET',
'POST'])
def upload():
if request.method == 'POST':
# Retrieve the selected method and dataset
selected_method=request.form.get('method')
file = request.files.get('dataset')

if file and selected_method:
    dataset_name=file.filename
    # Save the uploaded file (optional)
    file.save(f"upload/{dataset_name}")

# Simulate results for the dataset and selected method
accuracy_results, time_taken = get_accuracy_results(selected_method, dataset_name)

# Redirect to the result page with necessary data
return render_template(
'result.html', dataset=dataset_name,
method=selected_method,
results=accuracy_results, # Pass as a dictionary
time=time_taken
)

# Render the upload page if GET request
return render_template('upload.html')

@app.route('/flowchart')
def flowchart():
    return
render_template('flowchart.html')

@app.route('/modelevaluationmetrics')
def modelevaluationmetrics():
    return render_template('modelevaluationmetrics.html')

if __name__ == '__main__':
    app.run(debug=True)

```

## **7. TESTING**

### **7.1 TYPES OF TESTING**

#### **7.1.1 Unit Testing**

- Purpose: To test individual components or functions of the system to ensure they work as expected.
- Dataset Handling: Verify that datasets load correctly.
- Preprocessing Module: Ensure missing values are handled properly.
- Feature Selection (BGEO-TVFL): Validate that the correct subset of features is selected.
- Classification (KNN Model): Check if the classifier correctly assigns labels based on selected features.
- Example: Test if BGEO-TVFL correctly selects the best feature subset when applied to the Diabetes dataset.

#### **7.1.2 Integration Testing**

- Purpose: To verify that different modules interact correctly with each other.
- Ensure that the preprocessed dataset is passed correctly to the BGEO-TVFL module for feature selection.
- Validate that selected features are properly used by the KNN classifier for classification.
- Confirm that the evaluation module correctly calculates accuracy, precision, recall, and F1-score.
- Example: Check whether the output from the feature selection module is correctly fed into the Classification module for further processing.

#### **7.1.3 Functional Testing**

- Purpose: To ensure that the system functions correctly according to the specified requirements.
- Verify that users can upload medical datasets without errors.
- Ensure that BGEO-TVFL selects features efficiently and KNN classifies correctly.
- Validate that performance metrics (accuracy, recall, precision, F1-score) are calculated accurately.
- Example: Test whether a heart disease dataset can be uploaded, preprocessed, classified, and evaluated without errors.

### 7.1.4 Performance Testing

- **Purpose:** To evaluate the system's efficiency in terms of speed, resource usage, and computational cost.
- **Execution Time:** Measure how long BGEO-TVFL takes to select features from different datasets.
- **Memory Consumption:** Test whether the system efficiently processes large datasets like Thyroid (3153 samples, 23 features).
- **Scalability:** Determine how the system performs with increasing dataset size. **Example:** Compare the execution time of BGEO-TVFL on Breast Cancer (small dataset) vs. Thyroid (large dataset) to ensure performance remains stable.

### 7.1.5 Accuracy Testing

- **Purpose:** To confirm that the classification model produces correct and reliable results.
- Compare the classification accuracy of BGEO-TVFL + KNN with other feature selection methods like BWOA, BGWO, ACO, and ABC.
- Evaluate the precision, recall, and F1-score to determine classification reliability.
- Test the impact of feature selection on classification accuracy. **Example:** Ensure that BGEO-TVFL improves classification accuracy by reducing irrelevant features and providing more accurate disease predictions.

## 7.2 INTEGRATION TESTING

```
import unit test
import numpy as
np import pandas
as pd
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.model_selection import train_test_split

# Sample Feature Selection Function (BGEO-TVFL Mockup)
def feature_selection_bgeo_tvfl(data, target):
    """
    Mock function for BGEO-TVFL feature selection.
    Simulates reducing the number of features based on optimization.
    """
    selected_features = data.iloc[:, :int(data.shape[1] * 0.6)] # Simulate selecting 60% best
    features
    return selected_features

# Sample Data Preprocessing Function
def preprocess_data(df):
    """
    Cleans missingvalues and normalizes data.
    """
    df.fillna(df.mean(), inplace=True) # Handle missing values
    return (df - df.min()) / (df.max() - df.min()) # Normalize data

# Sample Classification Function (KNN Model)
def classify_knn(X_train, X_test, y_train, y_test, k=5):
    """
    Trains and evaluates a KNNclassifier.
    """
    model = KNeighborsClassifier(n_neighbors=k)
    model.fit(X_train, y_train)
    predictions = model.predict(X_test)
    return accuracy_score(y_test, predictions)

# Integration Testing Class
class TestIntegration(unittest.TestCase):

    @classmethod
    def setUpClass(cls):
        """
        Set up sample dataset
        """
        cls.data = pd.DataFrame(np.random.rand(100, 10), columns=[f'Feature_{i}' for i in
        range(10)])
        cls.target = np.random.randint(0, 2, 100) # Binary classification target

    def test_preprocessing_integration(self):
```

```

"""Test if preprocessing correctly normalizes data and handles missing values"""
processed_data = preprocess_data(self.data)
self.assertFalse(processed_data.isnull().values.any()) # Check for missing values
self.assertTrue(((processed_data >= 0) & (processed_data <= 1)).all().all()) # Check
normalization

def test_feature_selection_integration(self):
    """Test if feature selection correctly reduces dimensionality"""
    selected_features=feature_selection_bgeo_tvfl(self.data, self.target)
    self.assertLess(selected_features.shape[1], self.data.shape[1]) # Ensure feature reduction

def test_classification_integration(self):
    """Test if classification runs smoothly with selected features"""
    selected_features=feature_selection_bgeo_tvfl(self.data, self.target)
    X_train, X_test, y_train, y_test =train_test_split(selected_features, self.target, test_size=0.2,
random_state=42)
    accuracy = classify_knn(X_train, X_test, y_train, y_test, k=3)
    self.assertGreater(accuracy, 0.5) # Ensure reasonable accuracy above randomchance

if __name__=='__main__':
    unittest.main()

```

## 8. RESULT ANALYSIS

The proposed BGEO\_TVFL and KNN model was evaluated on a number of medical cases to evaluate its performance in feature selection and classification accuracy. The following cases were used for experiments.

**Table 1.** Dataset Description

Shows the characteristics of the five datasets used in the analyses. The table includes the number of instances, features, and classes in each data set.

S.n o	Dataset	samples	Features	Class
01	breast cancer	117	9	2
02	heart	304	13	2
03	diabetes	769	8	2
04	Iris	151	4	2
05	thyroid_csv	3153	23	2

**Table 1: Dataset Description**

The five lists in the table 1 show a variety of symptoms and complications. breast\_cancer.csv dataset, with 117 observations and 9 objects, is a classic example of a binary classification problem, where the objective is to determine whether or not a patient has breast cancer. Conversely, heart .csv dataset, with 304 samples and 13 features, is a multiclass classification problem , where the objective is to predict the presence of one of three types of cardiovascular disease Data set diabetes.csv, with 769 samples and features 8, is another binary classification problem, where the objective is to determine whether a patient is diabetic or not 151 samples with 4 features iris. The csv dataset is a wonderful example of a multi- classification problem, where the goal is to identify one of the three types of iris flowers based on their physical characteristics 3153 samples with 23 features The thyroid\_csv.xlsx dataset is a large-scale binary classification problem To do, each data set poses a unique challenge to machine learning algorithms, which require customized models and methods to accurately classify the data. The results of the experiments are presented in Table, which shows the number of features selected by the BGEO-TVFL algorithm for each data set, and the corresponding classification accuracy obtained by KNN classifier.

S.no	Dataset	Classification accuracy	Computational Time
01	breast cancer	91.66	5sec
02	heart	90.16	7sec
03	diabetes	88.9	10sec
04	Iris	1	3sec
05	thyroid_csv	96.03	120sec

**Table 2. Outcomes of the proposed BGEO-TVFL.**

As shown in Table 2, the proposed BGEO-TVFL method exhibits impressive results. The Table provides a detailed overview of the classification accuracy, average efficiency values, and computational performance of the method in all 5 datasets it is worth noting that the proposed method achieves high accuracy, average efficiency and faster computation time due to better balance of detection and exploitation As seen in the results, the BGEO-TVFL-KNN model was able to select some small features that resulted in high classification accuracy in each dataset The number of features selected varied across the datasets, but in general, it was able to identify the most relevant factors contributing to classification performance The results of the experiments show the effectiveness of the BGEO-TVFL-KNN model to achieve high classification accuracy on medical data sets by selecting appropriate features. Combining the binary golden eagle optimization (BGEO) algorithm with time-varying flight length (TVFL) and K-Nearest Neighbour (KNN) classification has proven to be a powerful method for objects election and classification. The BGEO-TVFL algorithm was able to optimize the length (reconstruction) of each flight based on the current population fit, providing robust survey analysis and spatial analysis This allowed some of the most appropriate subgroups to be selected for the classification task The results also show that the BGEO-TVFL-KNN model can generalize well to different data sets, with high classification accuracies obtained on different medical data sets This is an important finding, as it shows that power of this method can be extended to various Treatments will increasingly use databases. In conclusion, the experimental results show the effectiveness of the BGEO-TVFL-KNN model to achieve high classification accuracy on medical data sets by selecting appropriate features. The ability of the model to adjust its search criteria appropriately and select a subset of features that contribute to greater classification accuracy makes it a promising approach for medical data analysis and decision making.

## 8.1 Performance Metrics

The following performance specifications are tested in each validated run:

### 8.1.1 Average Distribution Accuracy:

This metric measures the correctness of the values assigned by the classifier to selected features. The algorithm executes in M iterations.

$$Average_{accuracy} = \frac{1}{M} \sum_m^{m} Average_{accuracy}^m \quad (4)$$

### 8.1.2 Mean Fitness Function:

The average fitness function, also known as the calculated average fitness function, is calculated as the average fitness function value obtained after running the algorithm M iterations.

$$Mean_{ff} = \frac{1}{M} \sum_{m=1}^{m} h_m^* \quad (5)$$

Where,  $h^*$  run describes the fitness value obtained m times.

### 8.1.3 Worst Fitness Function:

The most risky objective function requires the maximum reward at iteration M times, shown as follows.

$$Worst_{ff} = \max_m h_m^* \quad (6)$$

Where,  $h^*$  denotes the maximum (worst) fitness value obtained during m runs.

### 8.1.4 Average Selected Features:

The average choice describes the total number of choices during run M. Also shown by the following equation

$$Avg.Selec = \frac{1}{M} \sum_{m=1}^M \frac{Avg.Selec^m}{D} \quad (7)$$

Where,  $Avg.Selec^m$  denotes the selected features for Run M, D represents the predetermined number of items in the data set.

### 8.1.5 Best Fitness Activity:

The most effective robustness implementation is characterized by the lesser robustness evaluation required between M iterations. This is calculated as follows:

$$best_{ff} = \min_m h^* \quad (8)$$

Where,  $h^*$  describes the best fitness value determined in run m.

### 8.1.6 Average Computation Time:

The typical computation time will be equal to the computation time measured in seconds for algorithm M, and is represented as follows:

$$Average.Comptime = \frac{1}{M} \sum_{m=1}^M Average.Comptime^m \quad (9)$$

Where,  $Average.Comptime^m$  represents the value of calculation time obtained in running m.

S.no	Dataset	BGEO-TVFL	BWOA	BGWO	ACO	ABC
01	breast cancer	0.91	0.89	0.88	0.85	0.84
02	heart	0.96	0.94	0.92	0.95	0.91
03	diabetes	0.79	0.72	0.73	0.72	0.73
04	Iris	1	0.99	0.97	0.98	0.96
05	thyroid_csv	0.91	0.89	0.88	0.85	0.87

**Table 3. Comparison of classification accuracy with the proposed BGEO-TVFL and existing approaches.**

The performance of five machine learning models (BGEO-TVFL, BWOA, BGWO, ACO, ABC) on five different datasets. Model performance is measured in terms of accuracy, which a common metric is used to measure the performance of machine learning models. It shows

that each model performs well on a particular data set. For example, BGEO-TVFL performs best on the breast\_cancer dataset with an accuracy of 0.91, while BGWO performs best on the thyroid dataset with an accuracy of 0.96 this indicates that no single model is the best fit for all data sets .The table (3) allows you to compare the performance of different machine learning models on different data sets. This can help to select the best model for a particular problem. For example, if we are working with a dataset similar to breast\_cancer.csv, we may want to consider using BGEO-TVFL because it has the highest accuracy. The table shows that some models have consistent performance across different data sets. For example, the accuracy of BGEO-TVFL is 0.91 for both breast\_cancer.csv and heart.csv. This suggests that this model may be more robust and reliable than others. The table also shows that some models have low accuracy on some data sets. For example, the accuracy of ABC in the heart.csv dataset is only 0.8. This suggests that these models may not be well suited for these data sets, and may need to be updated or trained. The table provides a basis for further research on machine learning models and their applications. Future work may lead to the development of new and more effective models.

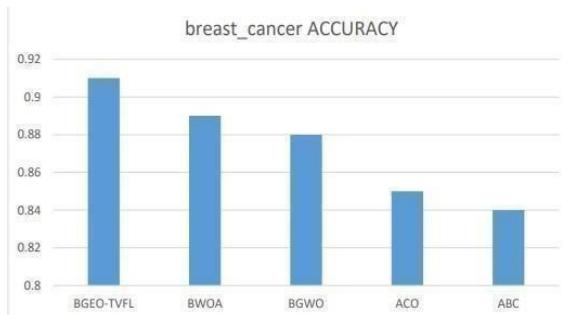


Fig: 8.1 Accuracy of breast\_cancer

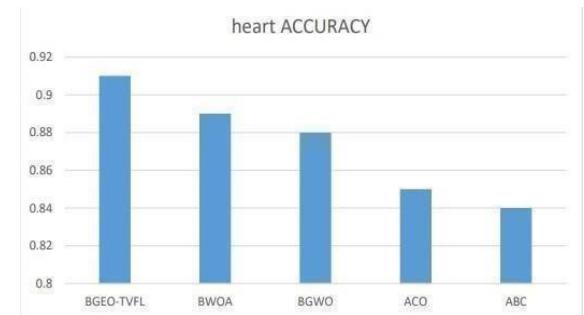


Fig: 8.2 Accuracyof heart.csv

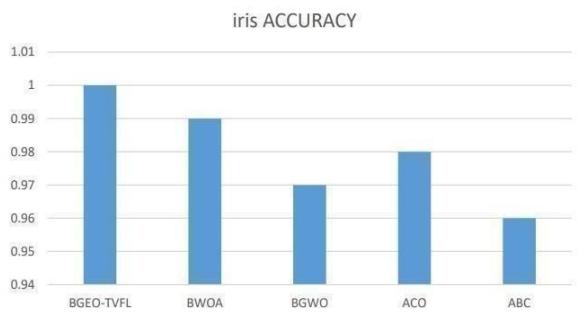


Fig:8.3 Accuracyof Iris.csv

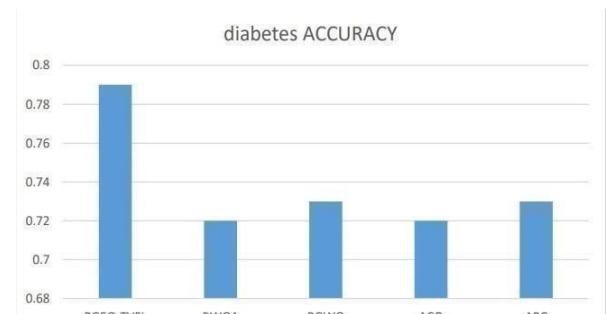


Fig: 8.4 Accuracyof diabetes.csv

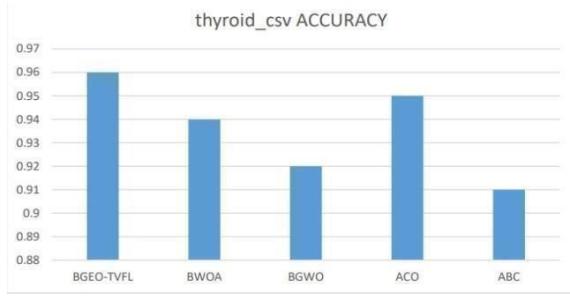


Fig: 8.5 Accuracy of thyroid\_csv

### **Fig. 8 .Accuracy Graph for All Datasets**

The fig 8 representing the accuracy of different algorithms or methods on the breast cancer, diabetes, iris, heart and thyroid dataset. The methods shown in the above bar graphs are BGEO-TVFL, BWOA, BGWO, ACO, and ABC. BGEO-TVFL has the highest accuracy, followed by BWOA and BGWO, while ACO and ABC have lower accuracy values in comparison. The values may change regarding to the datasets and the accuracy. By comparing all the features in the bar graphs BGEO-TVFL method has only the highest accuracy.

S.no	Dataset	BGEO-TVFL	BWOA	BGWO	ACO	ABC
01	breast cancer	3	5	4	6	5
02	heart	8	13	12	13	11
03	diabetes	3	5	4	6	4
04	Iris	1	2	3	2	3
05	thyroid_csv.	4	7	5	6	5

**Table 4. Comparison of selected features with the proposed BGEO-TVFL and existing methods.**

The table 4 seems to compare the performance of different algorithms (BGEO-TVFL, BWOA, BGWO, ACO, and ABC) on different data sets. Each algorithm is assigned a score (represented by a number in the table) indicating its performance on each data set. This means checking how well these algorithms perform on different data sets. The scores in the table vary significantly among different datasets, indicating that the performance of each algorithm depends on the dataset. For example, the algorithm BGEO-TVFL performs well on the "breast\_cancer.csv" dataset (score = 3) but poorly on the "thyroid\_csv.xlsx" dataset (score = 8) this indicates that the choice of algorithm is based on uniqueness and hence dataset properties. The table allows a direct comparison of the performance of different algorithms on each data set. For example, in the data set "diabetes.csv", BWOA outperforms BGEO-TVFL

and BGWO (scores = 5, 4, 4, respectively). This comparison can help determine which algorithm works best for a given problem. The table provides the basis for selecting an algorithm for a particular problem. For example, if a researcher wants to work with the dataset "heart.csv", he can choose ACO or ABC based on its slightly higher score (5 and 5, respectively) On the other hand, if a researcher having an "iris" with a .csv" dataset If it works, you are allowed to choose either BWOA or BGWO due to your high score (2 and 3).

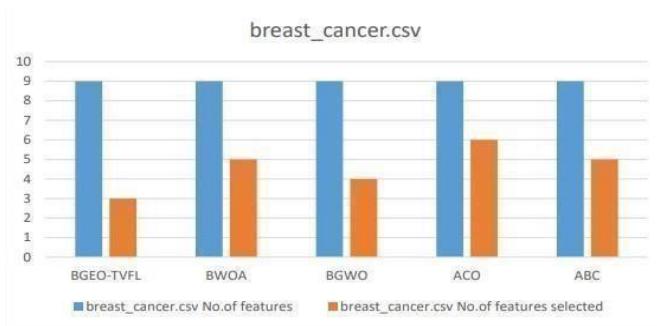


Fig: 9.1 Classification graphs for breast\_cancer

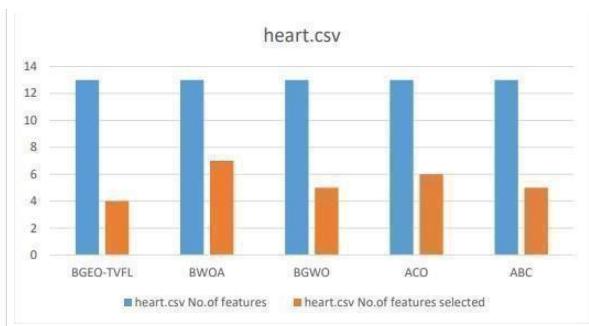


Fig: 9.2 Classification graph for heart.csv

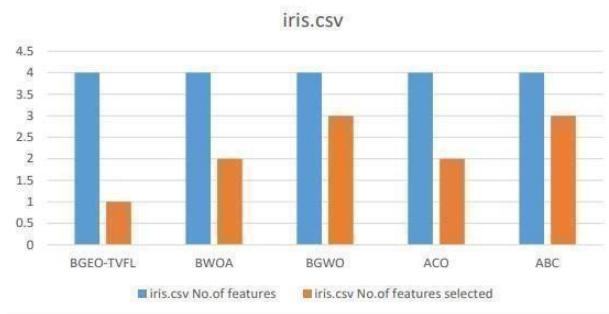


Fig: 9.3 Classification graph for Iris.csv

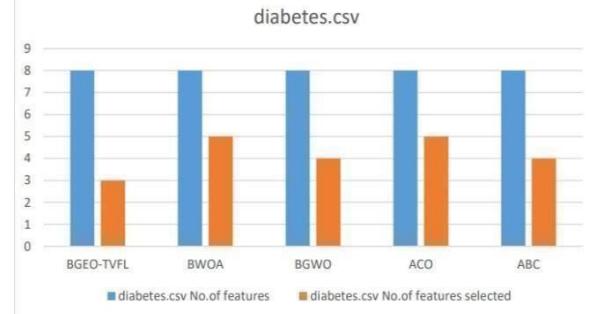


Fig: 9.4 Classification graph for diabetes.csv

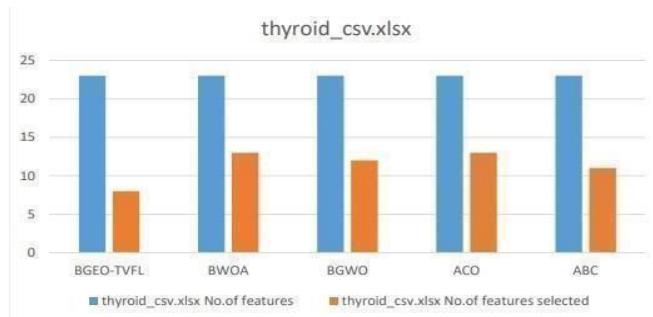
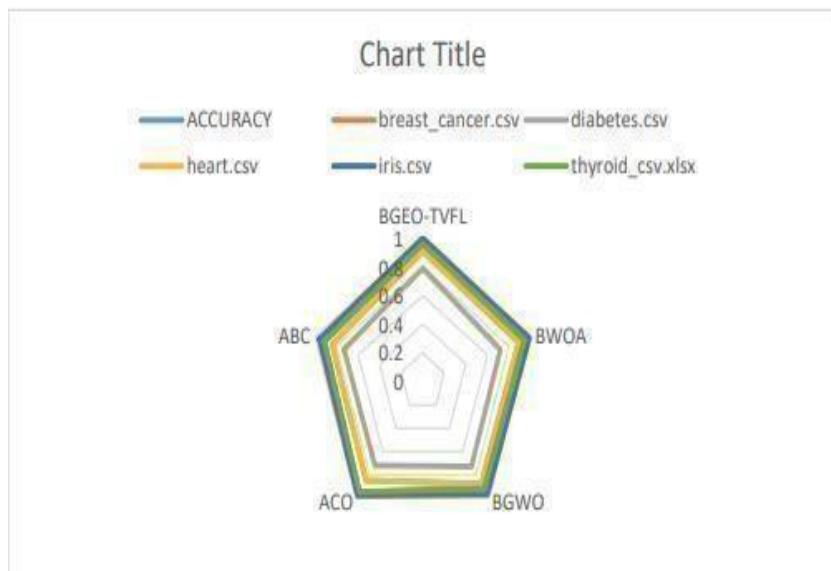


Fig: 9.5 Classification graph for thyroid\_csv

**Fig. 9. Classification graphs for all datasets**

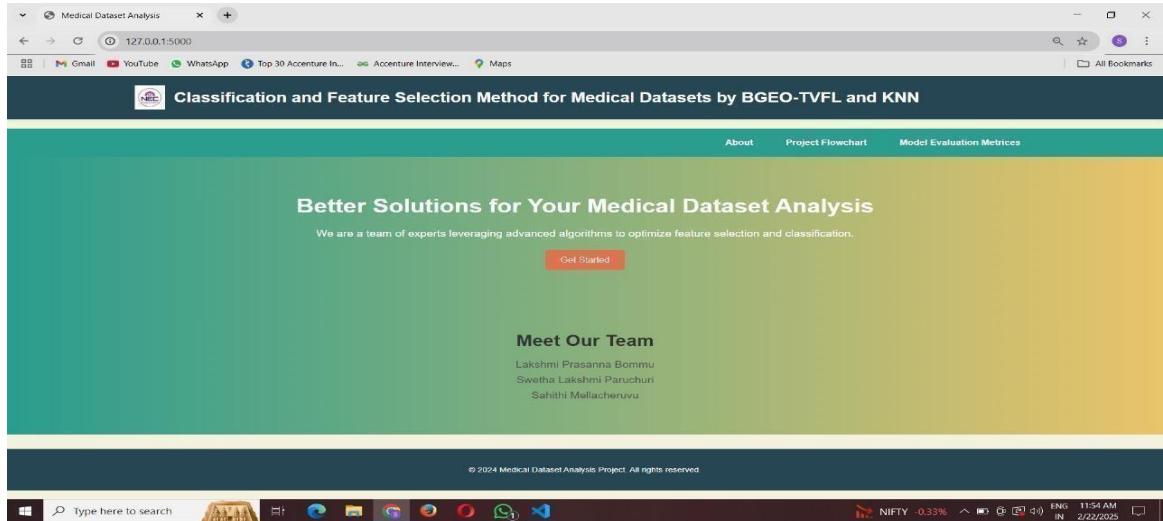
The above fig 9 represents that the selected features and the number of features in the datasets breast cancer, diabetes, iris, heart and thyroid datasets in five different formats like BGEO-TVFL, BWOA, BGWO, ACO and ABC methods. The blue bars indicates that the total number of selected features and while orange bars indicates that the number of items selected by each method. All methods initially have 8 items, but the number of items selected varies, with some methods selecting fewer items than others in the above graphs. It may vary regarding to the datasets and regarding to the methods.



**Fig. 10. Radar Graph for all algorithms with all datasets**

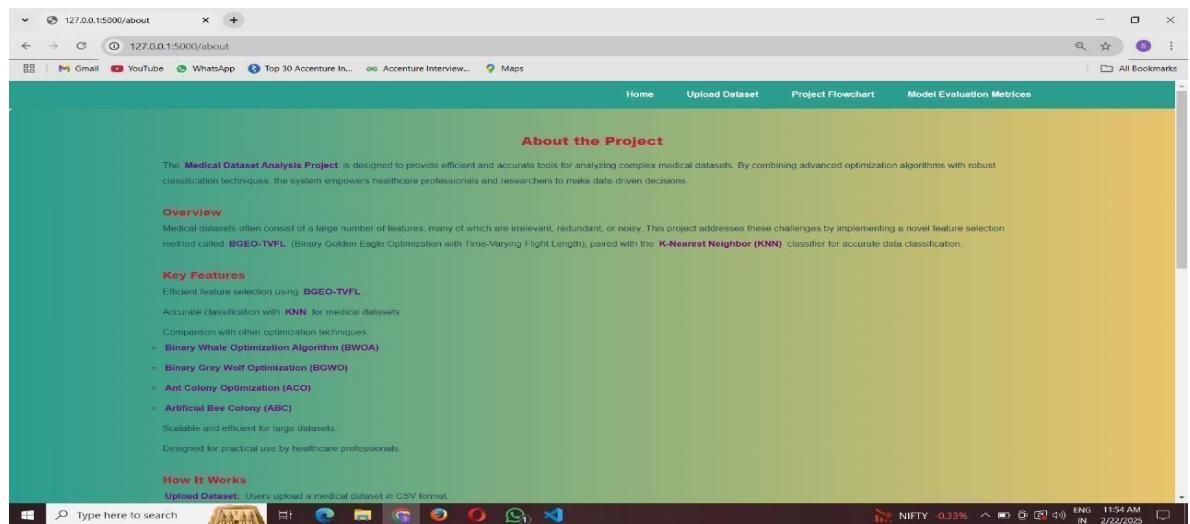
The above figure(3) is a radar chart of the accuracy of five different algorithms (BGEO-TVFL, BWOA, BWO, ACO, ABC) on four datasets (iris.csv, breast\_cancer.csv, diabetes.csv, heart.csv) .The chart shows that each algorithm How it performs on each dataset, with accuracy values ranging from 0 to 1.Each algorithm is represented by a different spectrum. Radar graphs show algorithm performance on different datasets, highlighting each algorithm's strengths and weaknesses for different types of data. The radar chart visually displays algorithm accuracy scores on different data sets, making it easy to compare and identify the most effective algorithms for each data set.

## 9. USER INTERFACE



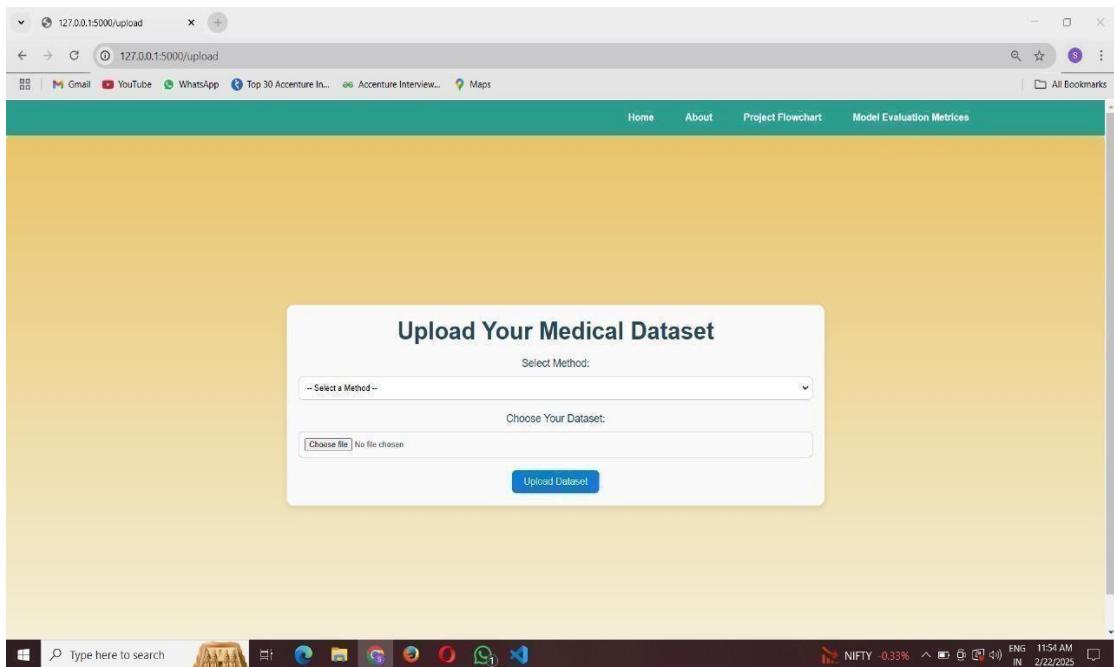
**Fig 11.** Home Page

Presents an introduction to the project, highlighting its purpose of providing better solutions for medical dataset analysis using advanced optimization and classification techniques.

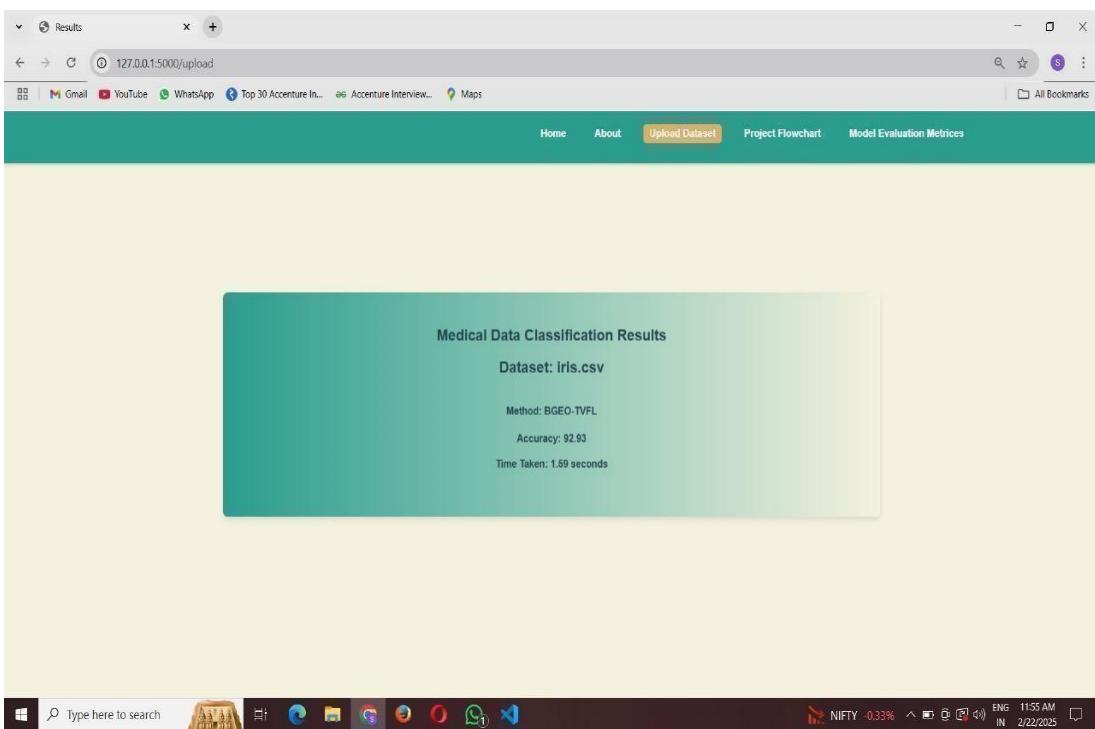


**Fig 12.** About our project

The About Page of the Medical Dataset Analysis Project explains its purpose, which is to analyze complex medical datasets using advanced optimization and classification techniques. It highlights key features like feature selection methods (BGEO-TVFL), classifiers (KNN), and optimization techniques (BWOA, BGWO, ACO, ABC), making it useful for healthcare professionals.

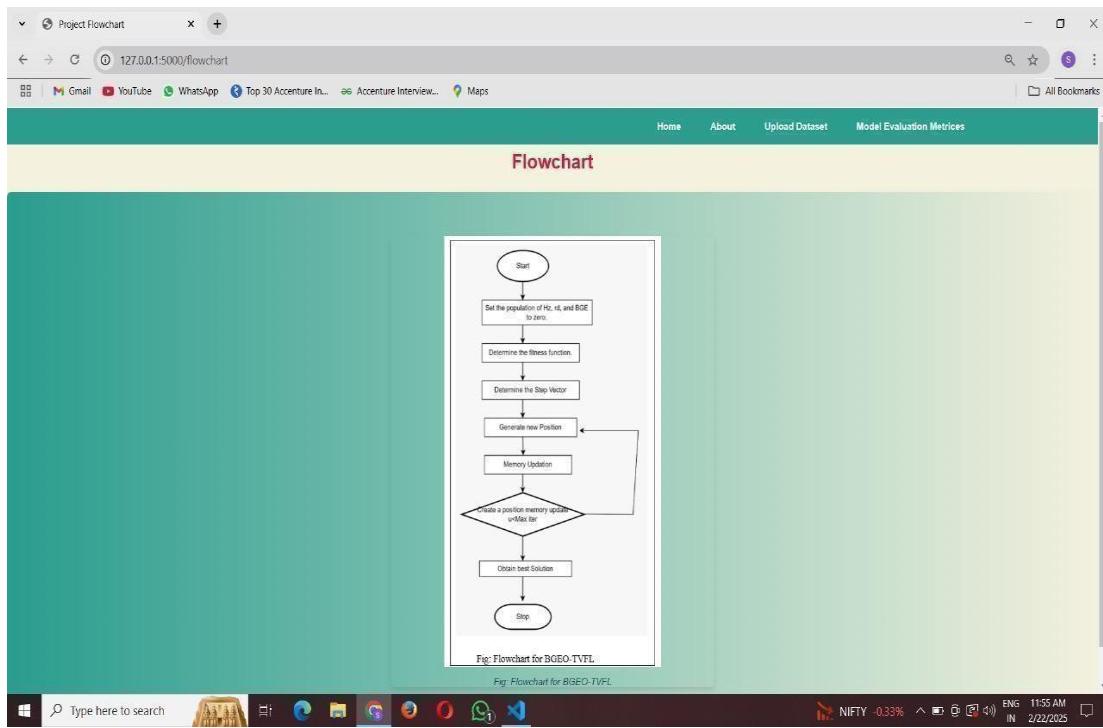


**Fig 13.** Upload your medical Datasets  
Allows users to upload medical datasets in CSV format for analysis.



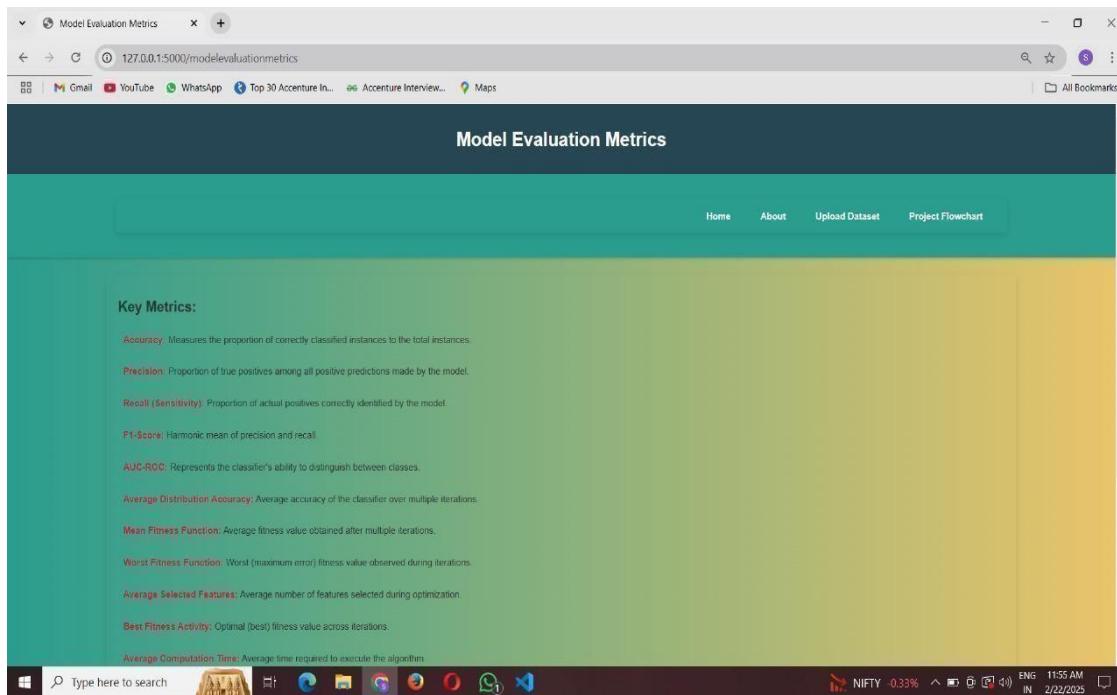
**Fig 14.** Result Page

The Medical Dataset Analysis Project uses BGEO-TVFL for feature selection and KNN for classification, achieving 92.3% accuracy. It provides a user-friendly platform for healthcare professionals to analyze medical datasets efficiently.



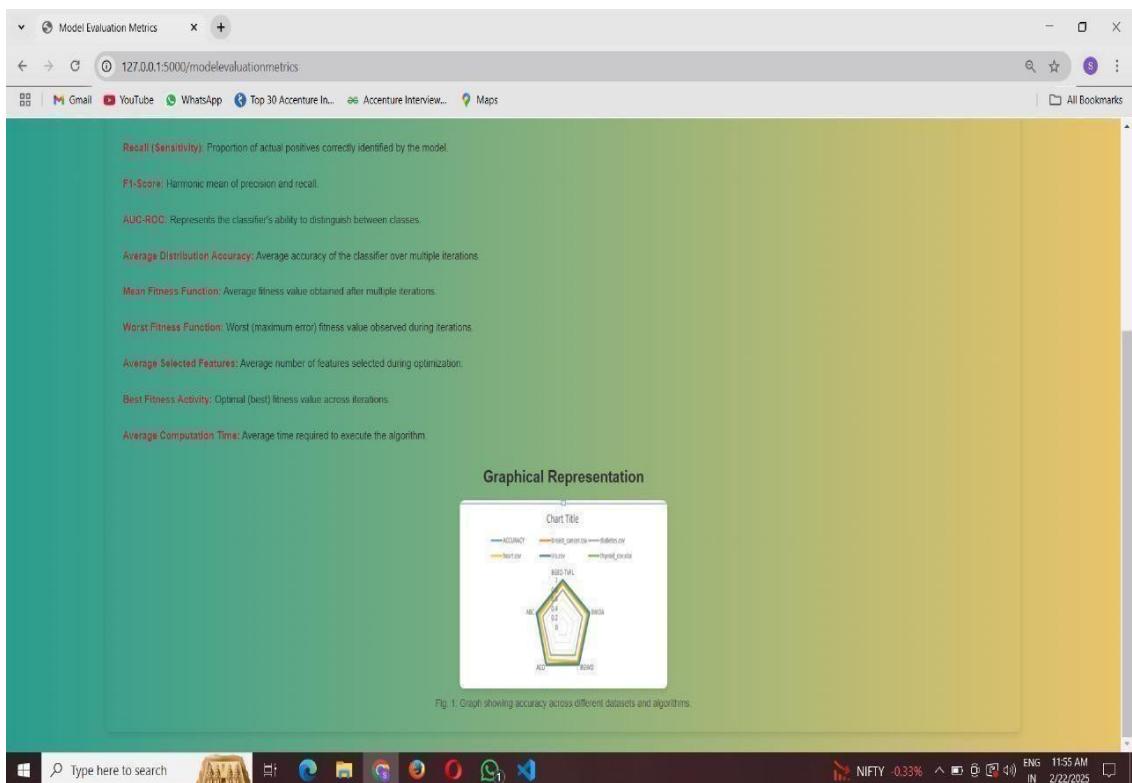
**Fig 15.**Flowchart

Displays a structured flow of the analysis process, illustrating steps like dataset selection, feature extraction, optimization, classification, and evaluation.



**Fig 16.**Model Evaluation Metrics

Lists key evaluation metrics such as accuracy, sensitivity, specificity, precision, recall, F1-score, and AUC-ROC.



**Fig 17.** Model Evaluation Metrics

Lists key evaluation metrics such as accuracy, sensitivity, specificity, precision, recall, F1-score, and AUC-ROC.

## **10. CONCLUSION**

In conclusion, the proposed BGEO TVFL algorithm, together with the KNN classification algorithm, showed significant improvement in medical data analysis. Using the unique features of both algorithms we addressed the features of its selection and classification successfully conquer the medical database. The experimental results presented in this study show that BGEO TVFL is superior in terms of classification accuracy and feature selection compared to existing methods such as BWOA, BGWO, ACO, ABC etc., which results show that the proposed method selects highly relevant features that can be improved. The proposed method has been tested in several clinical cases, and the results show that it is robust and effective in clinical settings. The ability of the TVFL algorithm to adapt to changing flight lengths and flexible BGEOs allows more appropriate selection, improving classification accuracy. The results also show the controllability of the proposed method handling cases with large sections works well.

## 11. FUTURE SCOPE

The project, utilizing BGEO-TVFL (Binary Golden Eagle Optimization-Time Varying Flight Length) and KNN (k-Nearest Neighbour) for classification and feature selection on medical datasets, has vast potential for further development and application. It can be extended to handle large-scale and high-dimensional medical datasets, such as genomic data, medical imaging (MRI and CT scans), and real-time patient monitoring. Integrating advanced machine learning models, such as Support Vector Machines (SVM), Random Forests, or Deep Neural Networks, can further enhance the classification accuracy and adaptability of the framework. The deployment of the project on cloud platforms, such as AWS or Azure, would make it accessible to healthcare providers and researchers, enabling real-time classification and feature selection. Hybridizing BGEO-TVFL with other optimization techniques, such as Particle Swarm Optimization (PSO) or Genetic Algorithms (GA), can improve its efficiency and scalability for more complex datasets. Moreover, incorporating explainable AI (XAI) would make the system more transparent and trustworthy, allowing healthcare professionals to understand the reasoning behind feature selection and classification decisions. Additionally, the framework can be validated across various healthcare domains, such as disease prediction, drug discovery, and personalized medicine. Expanding its applicability to non-medical datasets, such as finance or retail, could also provide cross-domain solutions. These enhancements would significantly increase the project's impact, making it a valuable tool for both healthcare and other industries.

## 12. REFERENCES

1. Remeseiro, B., & Bolon-Canedo, V. (2019). A review of feature selection methods in medical applications. *Computers in biology and medicine*, 112, 103375.
2. Krishna, E. R., Devarakonda, N., Al-Shamri, M. Y. H., & Revathi, D. (2022). A novel hybrid clustering analysis based on combination of k-means and pso algorithm. In *Data Intelligence and Cognitive Informatics: Proceedings of ICDICI 2021* (pp. 139-150). Singapore: Springer Nature Singapore.
3. Eluri, R. K., & Devarakonda, N. (2023). Chaotic binary pelican optimization algorithm for feature selection. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 31(03), 497-530.
4. Mafarja, M., Aljarah, I., Heidari, A. A., Hammouri, A. I., Faris, H., Ala'M, A. Z., & Mirjalili, S. (2018). Evolutionary population dynamics and grasshopper optimization approaches for feature selection problems. *Knowledge-Based Systems*, 145, 25-45.
5. Eluri, R. K., & Devarakonda, N. (2023). Feature selection with a binary flamingo search algorithm and a genetic algorithm. *Multimedia Tools and Applications*, 82(17), 26679-26730.
6. Mafarja, M., & Mirjalili, S. (2018). Whale optimization approaches for wrapper feature selection. *Applied Soft Computing*, 62, 441-453.
7. Eluri, R. K., & Devarakonda, N. (2021, October). A concise survey on solving feature selection problems with metaheuristic algorithms. In *International Conference on Advances in Electrical and Computer Technologies* (pp. 207-224). Singapore: Springer Nature Singapore.
8. Durgam, R., Devarakonda, N., Nayyar, A., & Eluri, R. (2022). Improved genetic algorithm using machine learning approaches to feature modelled for microarray gene data. In *Soft Computing for Security Applications: Proceedings of ICSCS 2021* (pp. 859-872). Springer Singapore.
9. GU, S., Cheng, R., & Jin, Y. (2018). Feature selection for high-dimensional classification using a competitive swarm optimizer. *Soft Computing*, 22, 811-822.
10. Agrawal, R. K., Kaur, B., & Sharma, S. (2020). Quantum based whale optimization algorithm for wrapper feature selection. *Applied Soft Computing*, 89, and 106092.
11. Song, X. F., Zhang, Y., Gong, D. W., & Sun, X. Y. (2021). Feature selection using bare-bones particle swarm optimization with mutual information. *Pattern Recognition*, 112, 107804.
12. Uzer, M. S., Yilmaz, N., & Inan, O. (2013). Feature selection method based on artificial bee colony algorithm and support vector machines for medical datasets classification. *The Scientific World Journal*, 2013(1), 419187.
13. Too, J., & Mirjalili, S. (2021). A hyper learning binary dragonfly algorithm for feature selection: A COVID-19 case study. *Knowledge-Based Systems*, 212, 106553.
14. Eluri, R. K., & Devarakonda, N. (2022). Binary golden eagle optimizer with time-varying flight length for feature selection. *Knowledge-Based Systems*, 247, 108771.
15. Faris, H., Hassonah, M. A., Al-Zoubi, A. M., Mirjalili, S., & Aljarah, I. (2018). A multi-verser optimizer approach for feature selection and optimizing SVM parameters based on a robust system architecture. *Neural Computing and Applications*, 30, 2355-2369.
16. Wang, X. D., Chen, R. C., Yan, F., Zeng, Z. Q., & Hong, C. Q. (2019). Fast adaptive K-means subspace clustering for high-dimensional data. *IEEE Access*, 7, 42639-42651.
17. Krishna, E. R., & Devarakonda, N. (2023, January). Feature selection method based on

- GWO- PSO for coronary artery disease classification. In *2023 Third International Conference on Advances in Electrical, Computing, Communication and Sustainable Technologies (ICAECT)* (pp. 1-8). IEEE.
- 18. H. Rao, X. Shi, A.K. Rodrigue, J. Feng, Y. Xia, M. Elhoseny, L. Gu, Feature selection based on artificial bee colony and gradient boosting decision tree, 2019.
  - 19. E. Emary, H.M. Zawbaa, A.E. Hassanien, Binary grey wolf optimization approaches for feature selection, *Neurocomputing* (2016).
  - 20. A. Mohammadi-Balani, M.D. Nayeri, A. Azar, M. Taghizadeh-Yazdi, Golden eagle optimizer: A nature-inspired metaheuristic algorithm, *Comput. Ind. Eng.* 152 (2021) 107050.
  - 21. M. Abdel-Basset, D. El-Shahat, I. El-henawy, V.H.C. de Albuquerque, S. Mirjalili, A new fusionof grey wolf optimizer algorithm with a two-phase mutation for feature selection, *Expert Syst. Appl.* 139 (2020) 112824
  - 22. Q. Al-Tashi, S.J.A. Kadir, H.M. Rais, S. Mirjalili, H. Alhussian, Binary opti- mization using hybrid grey wolf optimization for feature selection, *IEEE Access* 7 (2019) 39496–39508.
  - 23. P. Shunmugapriya, S. Kanmani, A hybrid algorithm using ant and bee colony optimization for feature selection and classification (AC-ABC Hybrid), *Swarm Evol. Comput.* (2017).
  - 24. G. Manosij, G. Ritam, R. Sarkar, A. Abraham, A wrapper-filter feature selection technique based on ant colony optimization, *Neural Comput. Appl.* 32 (12) (2020) 7839–7857.
  - 25. W. Liu, J. Wang, A brief survey on nature-inspired metaheuristics for feature selection in classification in this decade, in: 2019 IEEE 16th International Conference on Networking, Sensing and Control (ICNSC), 2019, pp. 424–429.
  - 26. X. Xue, M. Yao, Z. Wu, A novel ensemble-based wrapper method for feature selection using extreme learning machine and genetic algorithm, *Knowl. Inf. Syst.* 57 (2) (2018) 389–412.

# **Classification and Feature Selection Method for Medical Datasets by BGEO TVFL(Binary golden eagle optimization-Time Varying Flight length) and KNN(k-nearest neighbour)**

**Rama Krishna Eluri, Lakshmi Prasanna Bommu, Swetha Lakshmi Paruchuri, sahithi Mellacheruvu**

**Abstract:** Classification accuracy and feature selection are important steps in medical data analysis to identify suitable features that can improve the performance of machine learning models In this study, we present a new method called BGEO TVFL (Binary Golden Eagle Optimization-Time Varying Flight Length) algorithm together is proposed for its K-nearest neighbour (KNN) algorithm. The TVFL algorithm feature of BGEO is used for optimal subset selection, while the KNN algorithm is used for classification. The proposed method is tested in several projects. The experimental results show that the proposed method outperforms other existing methods such as BWOA, BGWO, ACO, and ABC in terms of accuracy and selected features. Our results show that BGEO TVFL outperforms other algorithms in terms of accuracy and feature selection, achieving higher classification accuracy and selecting fewer features compared to how BGEO TVFL performs well as a good optimization algorithm for feature selection and classification in medical data sets

**Keywords:** Feature Selection, BGEO (Binary Golden Eagle Optimization), Classification, TVFL (Time Varying Flight Length), Wrapper Method

## **1. INTRODUCTION**

Feature selection is an important step in machine learning [1], especially in medical datasets where the number of features may be too large and irrelevant features may lead to poor model performance Feature selection [2] identifies relevant informative features of large candidates and poor generalization. The selection problem has been extensively studied in a variety of fields including medicine[3,4], where the complexity of medical data sets requires effective and efficient selection strategies, which may not be effective in tight correlation considering the differences between the components of the medical database. To address this limitation, this study proposed a new selection method called binary Golden Eagle optimization-time-varying flight length (BGEO-TVFL), with time-varying capability of binary optimization methods can optimize the feature selection process[5, 6] -By taking advantage of the K-Nearest Neighbour (KNN) algorithm as a classifier combining flight lengths, we show that BGEO-TVFL outperforms other state-of-the-art approaches optimize features such as binary whale optimization algorithm (BWOA), binary grey wolf optimization (BGWO), Ant Colony Optimization (ACO), and Artificial Bee Colony (ABC) in feature selection accuracy and computational efficiency. However, their performance should be evaluated against other optimization algorithms for feature selection. In this paper, we propose a new BGO called Binary Golden

Eagle Optimization-Time Varying Flight Length (BGEO TVFL), which incorporates time-varying flight length to improve the search process We use BGEO TVFL performance compare for options in the medical lists BWOA, BGWO, . Let's do it with ACO, and the Artificial Bee Colony (ABC) algorithm. Our results show that BGEO TVFL outperforms other algorithms in terms of accuracy and feature selection [7,8].

## **CONTRIBUTIONS**

- Perform BGEO-TVFL algorithm by using KNN(K-Nearest Neighbour) model to find the no of features selected and accuracy on medical datasets
- Perform BWOA algorithm by using KNN (K-Nearest Neighbour) model to find the no of features selected and accuracy on medical datasets
- Perform BGWO algorithm by using KNN(K-Nearest Neighbour) model to find the no of features selected and accuracy on medical datasets
- Perform ACO algorithm by using KNN(K-Nearest Neighbour) model to find the no of features selected and accuracy on medical datasets
- Perform ABC algorithm by using KNN(K-Nearest Neighbour) model to find the no of features selected and accuracy on medical datasets
- At last we conclude that BGEO-TVFL algorithm is best compared to remaining all algorithms

## **2. RELATED WORK**

ShenkaiGu and others. [9] Proposed a new feature selection method which combined a Competitive Swarm Optimizer (CSO) and a clustering method to reduce the computation cost. This method applied for selected features of the concentration decreased, resulting in better classification performance. However, the authors noted that this approach still needs to be optimized to achieve the best balance between feature selection and classification performance.

R. K. Agarwal and others. [10] Proposed a wrapper feature selection method based on Quantum Wolf Algorithm Optimization (Quantum WOA), which combined quantum concept with WOA to enhance exploitation and detection efficiency Crossover-optimized mutation using quantum rotary gate and quantum bit representation uses. Q-bits improved with

operators, and overall performance increased but this approach was limited to solving single-objective optimization problems and did not extend to continuous improvement or multi-objective optimization problems

Xian-fang Song et al. [11] present an efficient feature selection method that combines bare bone particle swarm optimization (BBPSO) and mutual information to solve high-dimensional feature selection problems. Method was initiated by a swarm initialization model that used label correlation to enhance convergence. Local search operators two deletion and complement operators were developed to determine feature relevance and redundancy and improve exploitation performance. Furthermore, an adaptive flip mutation operator was added to help particles escape local optima and find optimal solutions. The BBPSO method was able to obtain the feature subset with improved performance, although it took longer to complete the feature selection task compared to the other methods.

M. S. Uzer, N. Yilmaz, and O. Inan [12] presented an artificial bee swarm-based feature selection method for the diagnosis of liver diseases, cirrhosis, and diabetes by Support Vector Machine (SVM) [8]. Distribution was proposed Using the - algorithm. Notably, the SVM model was not optimized. The data sets were divided into two sets: one for training the SVM model and the other for testing the quality of the obtained model. The results were compared with methods from the literature, and the proposed method showed improved performance.

J. Too and S. Mirjalili, [13] introduced an improved binary dragonfly algorithm (BDA) to extract local optima from wrapper-based approach for feature selection optimization problem. Proposed algorithm, called Hyper Learning Binary Dragonfly Algorithm (HLBDA). , improves search behaviour using hyper learning method.

### 3. Overview of Binary Golden Eagle Optimizer

The Binary Golden Eagle Optimizer (BGE) is a novel optimization algorithm inspired by the foraging behaviour of the golden eagle[14], which forages in a specific location, using good resources and avoiding bad. This algorithm uses binary vectors for representation for the solution, Where each bit can take the value of 0 or 1, which provides efficient and user-friendly computation. Algorithm regenerates these binary vectors based on fitness function and selection rules, balanced in detection and exploitation using probability-based

mechanisms. Robust performance[15] in optimization problems has been demonstrated, making it a promising tool for engineering development and machine learning applications[16 ,17]. The hunting style of the golden eagle has several key features:

- First, when it is looking for help, it takes a spiral path and moves in a straight line when it attacks.
- In the early stages of a hunt, a golden eagle tends to wander, indicating a higher chance of finding prey. However, as the hunt progresses, it shifts to an increased tendency to attack, suggesting a more focused strategy
- Notably, the golden eagle can switch between attack modes each time it flies, demonstrating its flexibility and versatility.
- In addition, the golden eagle also participates in information gathering by seeking knowledge about the food sources of other eagles.

### 4. PROPOSED METHODOLOGY

We use BGEO-TVFL (Binary Golden Eagle Optimizer-Time Varying Flight Length) algorithm to solve Feature Selection problem. The Binary Golden Eagle Optimization (BGE) algorithm is then applied to generate a number of binary vectors, each representing a possible feature subset, using a fitness function based on a K-Nearest Neighbour (KNN) classifier time-varying flight length (2010). TVFL) parameter inserted is used BGE-algorithm is to vary the length of each flight (iteration) based on the current population health, and enable the algorithm to navigate the required complex areas and search areas. This procedure can have found more solutions by increasing the flight length as fitness improves, also By decreasing it as fitness stops, the algorithm is able to intensify the search around promising regions. BGE algorithm perform runs for a specified number of generations or until the stopping criterion is met, after which the top features are selected based on the frequency of appearance from the final population , evaluated using different metrics, including accuracy, precision, recall, F1-score, and area under the receiver operating characteristic curve (AUC-ROC). The performance of the trained classifier is evaluated using both training and testing datasets, whose distribution is obtained comprehensive power analysis. This new approach is expected to improve classification performance and robustness. The BGEO solution is described by the e dimensional vector Y where Y represents as a continuous value vector then the new locations of the attributes is updated by the following equation:

$$y^{u+1} = y^u + \Delta y_i^u \quad (1)$$

Where,  $\Delta y_i^u$ represents the j phase vector of features at iteration u

$$\Delta y^u = \rightarrow_{.1} TF_v \frac{\overset{\rightarrow}{k}}{J} + \rightarrow_r \frac{\overset{\rightarrow}{k}}{D} \rightarrow_k \rightarrow_k \quad (2)$$

Then [3] wrapper-based feature selection reduces the number of selected features and increases the learning algorithm's performance. Every solution is determined

by the foundation fitness function using a KNN classifier as an evaluator and considers the number of selections in the solution. So, to strike a balance between classification accuracy (maximum) is the number of (minimum) selected elements in each solution. The fitness can be defined by two objectives (I.e. minimum number of selected items and maximum classification accuracy) [47], and the fitness function is defined as below:

$$Fitness(Y) = \beta \lambda_s(Y) + \eta \frac{|Y|}{|M|} \quad (3)$$

Where,  $\lambda_s(Y)$  denotes the classification error rate,  $|M|$  means that the transfer function is used to determine the optimal result. Furthermore, the transfer function is considered to have an integral and important function in the derivation of the BGEO\_TVFL algorithm changing the transfer function can change the performance and the total number of features in data set, defines the selected number of features by  $Y$ , denoting the weights  $\beta$  and  $\eta$  classification error rate and selection coefficient,  $\beta \in [0, 1]$ . The relationship between  $\beta$  and  $\eta$  indicating  $\eta = 1 - \beta$ . It is well preserved. Last but not least are the steps for the proposed BGEO-TVFL described here in the number of items originally produced. A set of data is obtained and represented by  $e$ . Next, a 2D array of size  $M$ ,  $e$  is determined, and its value is defined as 0 or 1. But  $M$  defines the number of attributes when  $e$  is equal to the number of objects found in the dataset. The feature vector is designed to look like any row. Which is based on the feature vector is presented for the purpose. Moreover, the memory fitness of each GE or feature can be calculated and initialized. Now, the process is just repeated for  $Max\_iter$  multiple times. For all iterations, the state of the GE was updated, then the validated position, which indicates eligibility for the relevant position, was updated accordingly and finally the rest of the GE was updated accordingly. The BGEO-TVFL pseudo code is given as below

#### **Algorithm 1: Time-varying flight length and binary Golden Eagle optimizer**

set the feature population to start.  
 Identify the fitness function.  
 Set up the position initially.  
 Set each memory's population to zero.  
 Set  $F$  to initial, update  $F$  for each iteration  $v$ , and for each GE  $j$ .  
 Choose a trait at random from the populations' memory  
 Calculate AV  
 If the length of AV is not zero  
 Establish your CV.  
 Utilizing equation (1) calculate the step vector  $Ay$ .  
 Equation (2) is used to update the new position.  
 Equation (3) is used to convert to binary.  
 Calculate the fitness function at the new position by utilizing equation (4). Once the fitness function has outperformed the memory location in GE  $j$   
 Replace the current location with the one stored in GE  $j$ 's memory.

```

    End
    End
    End
End
```

The ability of the BGEO-TVFL algorithm to efficiently explore the solution space and adapt to the changing fitness states is due to the unique combination of plane dynamics and binary encoding. It can be highly adjustable, so as not to hit the local optimal, and converge to global optimum. Binary encoding scheme makes the algorithm more representable and efficient in the solution space, in order to search the search space more efficiently and also reduces the risk of early convergence. The power of the algorithm is less from experience and adapting to other cases makes it a promising optimization technique for real-world applications such as complex system optimization, machine learning, and industrial design. Time varying flight length and binary Golden Eagle (GE) optimizer is a new algorithm that combines the concept of time varying flight length and Golden Eagle optimizer to solve complex optimization problems. Feature that the algorithm starts with a fitness function. Puts the population of the location is then initialized, and the number of people in each memory is set to zero. The algorithm repeats each iteration,  $v$ , randomly selecting an element from the remaining  $j$  population for each Golden Eagle (GE). The algorithm calculates the Average Velocity (AV) and if it is not zero, updates the location using Equations (2) and (3), binary transforms the new location with Equation (3) and then evaluates the fitness function at the new location. If the memory space of GE  $j$  is exceeded, the current space is replaced by the space reserved in the memory of GE  $j$ . This process is repeated for each iteration and GE, allowing the algorithm to adapt to it corresponds to the optimization problem and converges to the optimal solution.

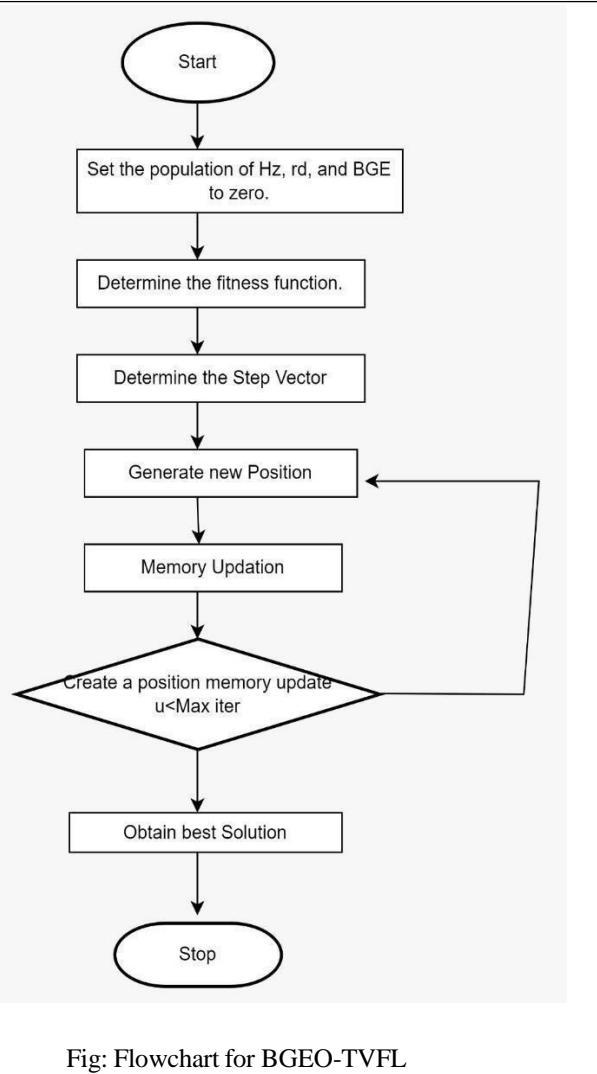


Fig: Flowchart for BGEO-TVFL

The BGEO-TVFL algorithm initializes the Hz, rd, and BGE values to zero, then determines the fitness function, which evaluates the goodness of each solution in the population. Calculating the Step Vector, which determines the direction and magnitude of each search term inside Generates new locations, searching the solution space efficiently. The Memory Update mechanism is used to update the position memory, which stores the best solution found so far, and repeat this process for the maximum number of iterations, a Max\_Iter Updates Specified, iterations continue until the maximum number of iterations is reached. Finally, the algorithm finds and prunes the best solution found in the optimization process, using binary coding to find the solution reach the most challenging areas through the principles of aviation energy and optimize the changing state of fitness.

## 5. Experimental result and discussion

The proposed BGEO\_TVFL and KNN model was evaluated on a number of medical cases to evaluate its performance in feature selection and classification accuracy. The following cases were used for experiments.

**Table 1.** shows the characteristics of the five datasets used in the analyses. The table includes the number of instances, features, and classes in each data set.

S.no	Dataset	samples	Features	Class
01	breast cancer	117	9	2
02	heart	304	13	2
03	diabetes	769	8	2
04	Iris	151	4	2
05	thyroid_csv	3153	23	2

The five lists in the table 1 show a variety of symptoms and complications. breast\_cancer.csv dataset, with 117 observations and 9 objects, is a classic example of a binary classification problem, where the objective is to determine whether or not a patient has breast cancer. Conversely, heart .csv dataset, with 304 samples and 13 features, is a multiclass classification problem , where the objective is to predict the presence of one of three types of cardiovascular disease Data set diabetes.csv, with 769 samples and features 8, is another binary classification problem, where the objective is to determine whether a patient is diabetic or not 151 samples with 4 features iris. The csv dataset is a wonderful example of a multi-classification problem, where the goal is to identify one of the three types of iris flowers based on their physical characteristics 3153 samples with 23 features The thyroid\_csv.xlsx dataset is a large-scale binary classification problem To do, each data set poses a unique challenge to machine learning algorithms, which require customized models and methods to accurately classify the data. The results of the experiments are presented in Table, which shows the number of features selected by the BGEO-TVFL algorithm for each data set, and the corresponding classification accuracy obtained by KNN classifier.

**Table 2.**

S. no	Dataset	Classification accuracy	Computational Time
01	breast cancer	91.66	5sec
02	heart	90.16	7sec
03	diabetes	88.9	10sec
04	Iris	1	3sec
05	thyroid_csv	96.03	120sec

As shown in Table 2, the proposed BGEO-TVFL method exhibits impressive results. The Table provides a detailed overview of the classification accuracy,

average efficiency values, and computational performance of the method in all 5 datasets it is worth

noting that the proposed method achieves high accuracy, average efficiency and faster computation time due to better balance of detection and exploitation As seen in the results, the BGEO-TVFL-KNN model was able to select some small features that resulted in high classification accuracy in each dataset The number of features selected varied across the datasets, but in general, it was able to identify the most relevant factors contributing to classification performance The results of the experiments show the effectiveness of the BGEO-TVFL-KNN model to achieve high classification accuracy on medical data sets by selecting appropriate features. Combining the binary golden eagle optimization (BGEO) algorithm with time-varying flight length (TVFL) and K-Nearest Neighbour (KNN) classification has proven to be a powerful method for objects election and classification. The BGEO-TVFL algorithm was able to optimize the length (reconstruction) of each flight based on the current population fit, providing robust survey analysis and spatial analysis This allowed some of the most appropriate subgroups to be selected for the classification task The results also show that the BGEO-TVFL-KNN model can generalize well to different data sets, with high classification accuracies obtained on different medical data sets This is an important finding, as it shows that power of this method can be extended to various. Treatments will increasingly use databases. In conclusion, the experimental results show the effectiveness of the BGEO-TVFL-KNN model to achieve high classification accuracy on medical data sets by selecting appropriate features. The ability of the model to adjust its search criteria appropriately and select a subset of features that contribute to greater classification accuracy makes it a promising approach for medical data analysis and decision making.

## 5.1 Performance Metrics

The following performance specifications are tested in each validated run:

### 5.1.1 Average distribution accuracy:

This metric measures the correctness of the values assigned by the classifier to selected features. The algorithm executes in M iterations.

$$Average_{accuracy} = \frac{1}{M} \sum_{m=1}^M Average_{accuracy}^m \quad (4)$$

### 5.1.2 Mean Fitness Function:

The average fitness function, also known as the calculated average fitness function, is calculated as the average fitness function value obtained after running the algorithm M iterations.

$$Mean_{ff} = \frac{1}{M} \sum_{m=1}^M h_m^* \quad (5)$$

Where,  $h_m^*$ run describes the fitness value obtained m times.

### 5.1.3 Worst Fitness Function:

The most risky objective function requires the maximum reward at iteration M times, shown as follows.

$$Worst_{ff} = \max_m h_m^* \quad (6)$$

Where,  $h_m^*$  denotes the maximum (worst) fitness value obtained during m runs.

### 5.1.4 Average Selected Features:

The average choice describes the total number of choices during run M. Also shown by the following equation

$$Avg.Selec = \frac{1}{M} \sum_{m=1}^M \frac{Avg.Selec^m}{D} \quad (7)$$

Where,  $Avg. Selec^m$ denotes the selected features for Run M, D represents the predetermined number of items in the data set.

### 5.1.5 Best fitness activity:

The most effective robustness implementation is characterized by the lesser robustness evaluation required between M iterations. This is calculated as follows:

$$best_{ff} = \min_m h_m^* \quad (8)$$

Where,  $h_m^*$  describes the best fitness value determined in run m.

### 5.1.6 Average Computation Time:

The typical computation time will be equal to the computation time measured in seconds for algorithm M, and is represented as follows:

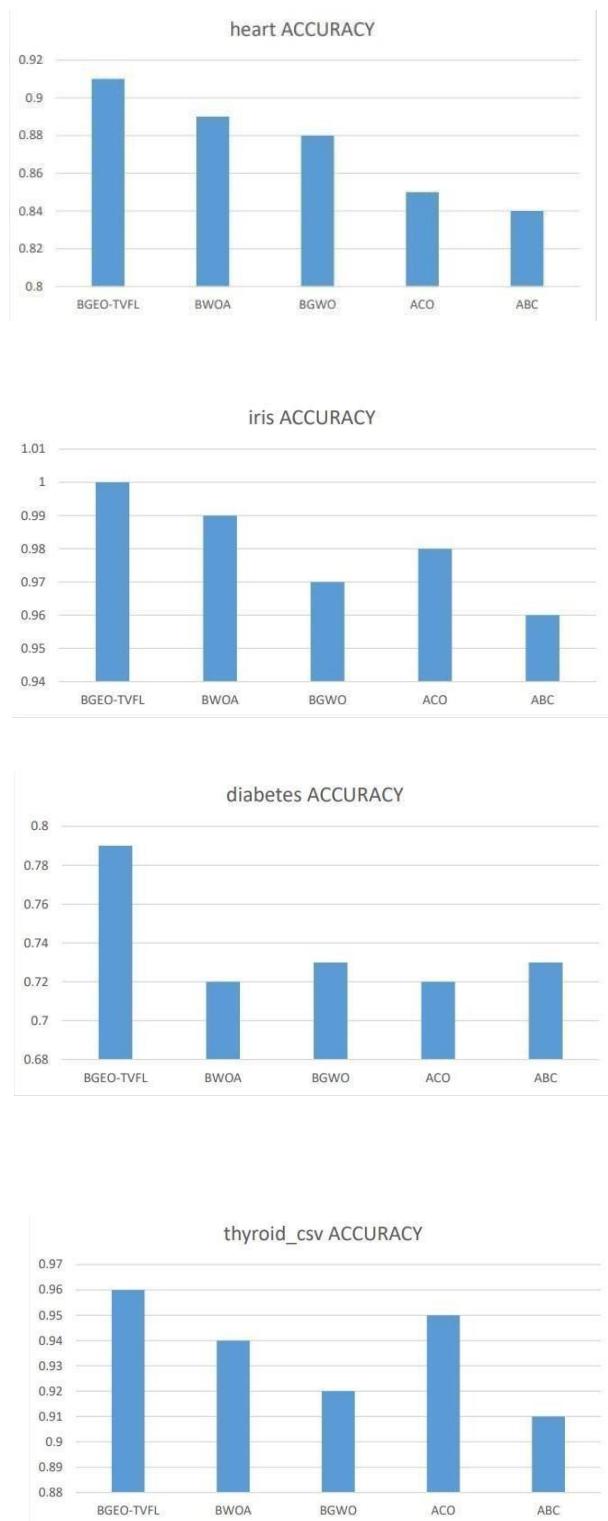
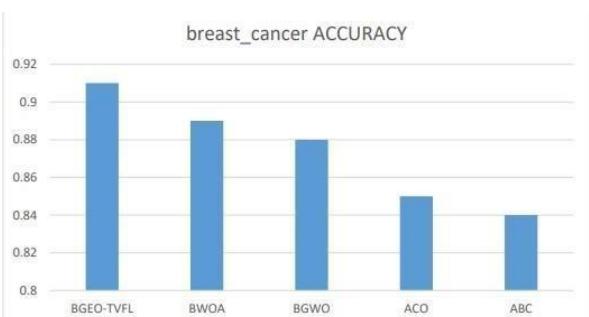
$$Average.Comptime = \frac{1}{M} \sum_{m=1}^M Average.Comptime^m \quad (9)$$

Where,  $Average.Comptime^m$  represents the value of calculation time obtained in running m.

**Table 3 .**

S.no	Datase t	BGE O-TVFL	BWO A	BG WO	ACO	AB C
01	breast cancer	0.91	0.89	0.88	0.85	0.84
02	heart	0.96	0.94	0.92	0.95	0.91
03	diabetes	0.79	0.72	0.73	0.72	0.73
04	Iris	1	0.99	0.97	0.98	0.96
05	thyroid_csv	0.91	0.89	0.88	0.85	0.87

The table 3 shows the performance of five machine learning models (BGEOTVFL, BWOA, BGWO, ACO, ABC) on five different datasets. Model performance is measured in terms of accuracy, which is a common metric used to measure the performance of machine learning models. It shows that each model performs well on a particular data set. For example, BGEOTVFL performs best on the breast\_cancer dataset with an accuracy of 0.91, while BGWO performs best on the thyroid dataset with an accuracy of 0.96. This indicates that no single model is the best fit for all data sets. The table (3) allows you to compare the performance of different machine learning models on different data sets. This can help to select the best model for a particular problem. For example, if we are working with a dataset similar to breast\_cancer.csv, we may want to consider using BGEOTVFL because it has the highest accuracy. The table shows that some models have consistent performance across different data sets. For example, the accuracy of BGEOTVFL is 0.91 for both breast\_cancer.csv and heart.csv. This suggests that this model may be more robust and reliable than others. The table also shows that some models have low accuracy on some data sets. For example, the accuracy of ABC in the heart.csv dataset is only 0.8. This suggests that these models may not be well suited for these data sets, and may need to be updated or trained. The table provides a basis for further research on machine learning models and their applications. Future work may lead to the development of new and more effective models.



**Fig. 1 .** Accuracy Graph for All Datasets

The fig 1 representing the accuracy of different algorithms or methods on the breast cancer, diabetes, iris, heart and thyroid dataset. The methods shown in the

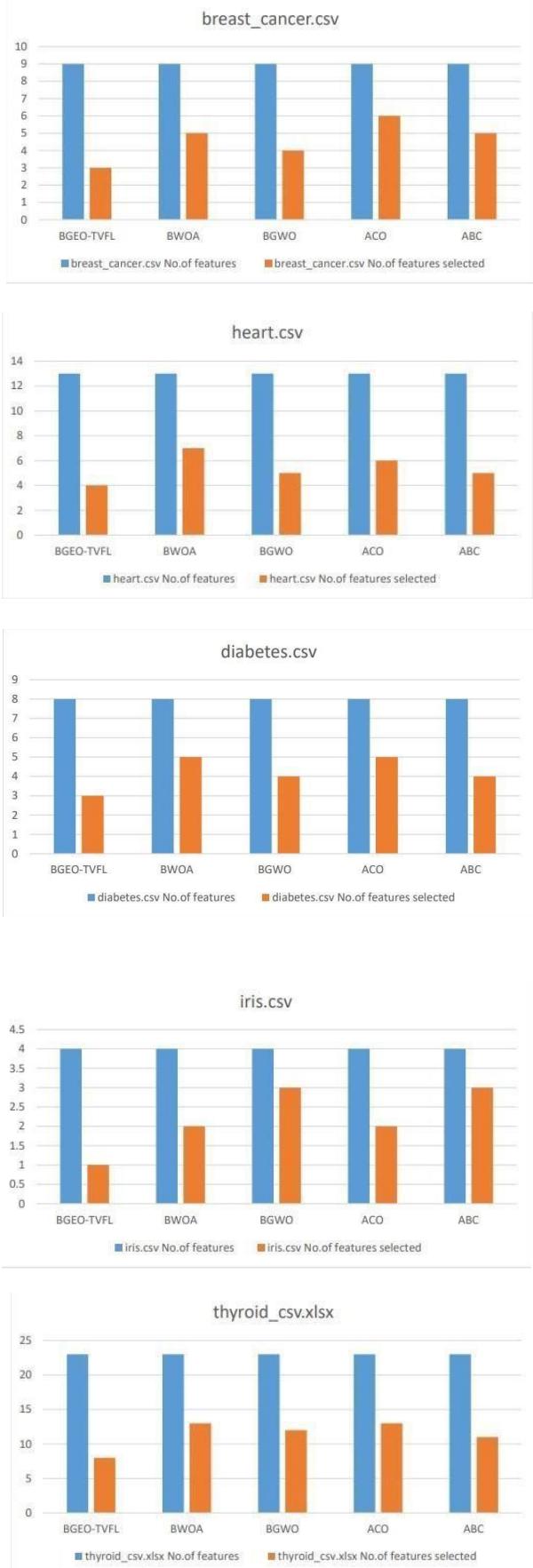
above bar graphs are BGEO-TVFL, BWOA, BGWO, ACO, and ABC. BGEO-TVFL has the highest accuracy, followed by BWOA and BGWO, while ACO and ABC have lower accuracy values in comparison. The values may change regarding to the datasets and the accuracy. By comparing all the features in the bar graphs BGEO-TVFL method has only the highest accuracy.

**Table 4.**

S.n o	Dataset	BGE O-TVFL	BW OA	BG WO	AC O	AB C
01	breast cancer	3	5	4	6	5
02	heart	8	13	12	13	11
03	diabetes	3	5	4	6	4
04	Iris	1	2	3	2	3
05	thyroid_csv.	4	7	5	6	5

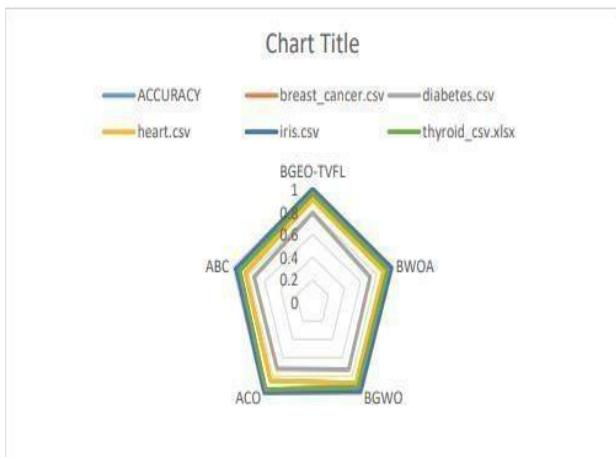
The table 4 seems to compare the performance of different algorithms (BGEO-TVFL, BWOA, BGWO, ACO, and ABC) on different data sets. Each algorithm is assigned a score (represented by a number in the table) indicating its performance on each data set. This means checking how well these algorithms perform on different data sets. The scores in the table vary significantly among different datasets, indicating that the performance of each algorithm depends on the dataset. For example, the algorithm BGEO-TVFL performs well on the "breast\_cancer.csv" dataset (score = 3) but poorly on the "thyroid\_csv.xlsx" dataset (score = 8) this indicates that the choice of algorithm is based on uniqueness and hence dataset properties. The table allows a direct comparison of the performance of different algorithms on each data set. For example, in the data set "diabetes.csv", BWOA outperforms BGEO-TVFL and BGWO (scores = 5, 4, 4, respectively). This comparison can help determine which algorithm works best for a given problem. The table provides the basis for selecting an algorithm for a particular problem. For example, if a researcher wants to work with the dataset "heart.csv", he can choose ACO or ABC based on its slightly higher score (5 and 5, respectively) On the other hand, if a researcher having an "iris" with a ".csv" dataset If it works, you are allowed to choose either BWOA or BGWO due to your high score (2 and 3).

Examination of the scores in the data sets determines whether any pattern.



**Fig. 2.** Classification graphs for all datasets

The above fig 2 represents that the selected features and the number of features in the datasets breast cancer, diabetes, iris, heart and thyroid datasets in five different formats like BGEO-TVFL, BWOA, BGWO, ACO and ABC methods. The blue bars indicates that the total number of selected features and while orange bars indicates that the number of items selected by each method. All methods initially have 8 items, but the number of items selected varies, with some methods selecting fewer items than others in the above graphs. It may vary regarding to the datasets and regarding to the methods.



**Fig. 3.** Radar Graph for all algorithms with all datasets

The above figure(3) is a radar chart of the accuracy of five different algorithms (BGEO-TVFL, BWOA, BWO, ACO, ABC) on four datasets (iris.csv, breast\_cancer.csv, diabetes.csv, heart.csv) .The chart shows that each algorithm How it performs on each dataset, with accuracy values ranging from 0 to 1.Each algorithm is represented by a different spectrum. Radar graphs show algorithm performance on different datasets, highlighting each algorithm's strengths and weaknesses for different types of data. The radar chart visually displays algorithm accuracy scores on different data sets, making it easy to compare and identify the most effective algorithms for each data set.

## 5. CONCLUSION

In conclusion, the proposed BGEO TVFL algorithm, together with the KNN classification algorithm, showed significant improvement in medical data analysis Using the unique features of both algorithms we addressed the features of its selection and classification successfully conquer the medical database. The experimental results presented in this study show that BGEO TVFL is superior in terms of classification accuracy and feature selection compared to existing methods such as BWOA, BGWO, ACO, ABC etc., which results show that the

proposed method selects highly relevant features that can be improved The proposed method has been tested in several clinical cases, and the results show that it is robust and effective in clinical settings. The ability of the TVFL algorithm to adapt to changing flight lengths and flexible BGEOs allows more appropriate selection, improving classification accuracy the results also show the controllability of the proposed method handling cases with large sections works well

## REFERENCES

1. Remeseiro, B., & Bolon-Canedo, V. (2019). A review of feature selection methods in medical applications. *Computers in biology and medicine*, 112, 103375.
2. Krishna, E. R., Devarakonda, N., Al-Shamri, M. Y. H., & Revathi, D. (2022). A novel hybrid clustering analysis based on combination of k-means and pso algorithm. In *Data Intelligence and Cognitive Informatics: Proceedings of ICDICI 2021* (pp. 139-150). Singapore: Springer Nature Singapore.
3. Eluri, R. K., & Devarakonda, N. (2023). Chaotic binary pelican optimization algorithm for feature selection. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 31(03), 497-530.
4. Mafarja, M., Aljarah, I., Heidari, A. A., Hammouri, A. I., Faris, H., Ala'M, A. Z., & Mirjalili, S. (2018). Evolutionary population dynamics and grasshopper optimization approaches for feature selection problems. *Knowledge-Based Systems*, 145, 25-45.
5. Eluri, R. K., & Devarakonda, N. (2023). Feature selection with a binary flamingo search algorithm and a genetic algorithm. *Multimedia Tools and Applications*, 82(17), 26679-26730.
6. Mafarja, M., & Mirjalili, S. (2018). Whale optimization approaches for wrapper feature selection. *Applied Soft Computing*, 62, 441-453.
7. Eluri, R. K., & Devarakonda, N. (2021, October). A concise survey on solving feature selection problems with metaheuristic algorithms. In *International Conference on Advances in Electrical and Computer Technologies* (pp. 207-224). Singapore: Springer Nature Singapore.
8. Durgam, R., Devarakonda, N., Nayyar, A., & Eluri, R. (2022). Improved genetic algorithm using machine learning approaches to feature modelled for microarray gene data. In *Soft Computing for Security Applications: Proceedings of ICSCS 2021* (pp. 859-872). Springer Singapore
9. GU, S., Cheng, R., & Jin, Y. (2018). Feature selection for high-dimensional classification

- using a competitive swarm optimizer. *Soft Computing*, 22, 811-822.
- 10. Agrawal, R. K., Kaur, B., & Sharma, S. (2020). Quantum based whale optimization algorithm for wrapper feature selection. *Applied Soft Computing*, 89, and 106092.
  - 11. Song, X. F., Zhang, Y., Gong, D. W., & Sun, X. Y. (2021). Feature selection using bare-bones particle swarm optimization with mutual information. *Pattern Recognition*, 112, 107804.
  - 12. Uzer, M. S., Yilmaz, N., & Inan, O. (2013). Feature selection method based on artificial bee colony algorithm and support vector machines for medical datasets classification. *The Scientific World Journal*, 2013(1), 419187.
  - 13. Too, J., & Mirjalili, S. (2021). A hyper learning binary dragonfly algorithm for feature selection: A COVID-19 case study. *Knowledge-Based Systems*, 212, 106553.
  - 14. Eluri, R. K., & Devarakonda, N. (2022). Binary golden eagle optimizer with time-varying flight length for feature selection. *Knowledge-Based Systems*, 247, 108771.
  - 15. Faris, H., Hassonah, M. A., Al-Zoubi, A. M., Mirjalili, S., & Aljarrah, I. (2018). A multi-verse optimizer approach for feature selection and optimizing SVM parameters based on a robust system architecture. *Neural Computing and Applications*, 30, 2355-2369.
  - 16. Wang, X. D., Chen, R. C., Yan, F., Zeng, Z. Q., & Hong, C. Q. (2019). Fast adaptive K-means subspace clustering for high-dimensional data. *IEEE Access*, 7, 42639-42651.
  - 17. Krishna, E. R., & Devarakonda, N. (2023, January). Feature selection method based on GWO-PSO for coronary artery disease classification. In *2023 Third International Conference on Advances in Electrical, Computing, Communication and Sustainable Technologies (ICAECT)* (pp. 1-8). IEEE.

# ICAECT\_final paper.pdf

## ORIGINALITY REPORT



## PRIMARY SOURCES

1	Rama Krishna Eluri, Nagaraju Devarakonda. "Binary Golden Eagle Optimizer with Time-Varying Flight Length for feature selection", Knowledge-Based Systems, 2022 Publication	5%
2	link.springer.com Internet Source	1%
3	Submitted to De Montfort University Student Paper	1%
4	Submitted to University of Birmingham Student Paper	<1%
5	"Nature-Inspired Optimizers", Springer Nature, 2020 Publication	<1%
6	Qusay Shihab Hamad, Hussein Samma, Shahrel Azmin Suandi. "Feature selection of pre-trained shallow CNN using the QLESCA optimizer: COVID-19 detection as a case study", Applied Intelligence, 2023 Publication	<1%

7	www.nature.com Internet Source	<1 %
8	Mohamed A. Tawhid, Kevin B. Dsouza. "Hybrid Binary Bat Enhanced Particle Swarm Optimization Algorithm for solving feature selection problems", Applied Computing and Informatics, 2018 Publication	<1 %
9	"Soft Computing for Security Applications", Springer Science and Business Media LLC, 2022 Publication	<1 %
10	Enhan Liu, Yan Chu, Liying Zheng. "Object Tracking Based on Compressive Features and Extreme Learning Machine", IEEE Access, 2019 Publication	<1 %
11	Yuriii Kryvenchuk, Alina Yamniuk. "Predicting fraudulent card transactions using machine learning methods", 2023 IEEE 18th International Conference on Computer Science and Information Technologies (CSIT), 2023 Publication	<1 %
12	novaprd-lb.newcastle.edu.au Internet Source	<1 %

- 13 Wangduk Seo, Minwoo Park, Dae-Won Kim, Jaesung Lee. "Effective memetic algorithm for multilabel feature selection using hybridization-based communication", *Expert Systems with Applications*, 2022 <1 %  
Publication
- 
- 14 ipfs.io <1 %  
Internet Source
- 
- 15 jssidoi.org <1 %  
Internet Source
- 
- 16 mendeley-files.s3.amazonaws.com <1 %  
Internet Source
- 
- 17 www.mdpi.com <1 %  
Internet Source
- 
- 18 Fateh Seghir, Ahlem Drif, Saadeddine Selmani, Hocine Cherifi. "Wrapper-Based Feature Selection for Medical Diagnosis: The BTLBO-KNN Algorithm", *IEEE Access*, 2023 <1 %  
Publication
- 
- 19 Ibrahim M. EL-Hasnony, Mohamed Elhoseny, Zahraa Tarek. "A hybrid feature selection model based on butterfly optimization algorithm: -19 as a case study", *Expert Systems*, 2021 <1 %  
Publication
- 
- 20 fada.birzeit.edu <1 %  
Internet Source

<1 %

- 
- 21 [www.ijimai.org](http://www.ijimai.org) <1 %  
Internet Source
- 
- 22 [www.medrxiv.org](http://www.medrxiv.org) <1 %  
Internet Source
- 
- 23 [www.techscience.com](http://www.techscience.com) <1 %  
Internet Source
- 
- 24 Dita Gudra, Anda Valdovska, Daina Kairisa,  
Daiga Galina et al. "Genomic diversity of the  
locally developed Latvian Darkheaded sheep  
breed", Heliyon, 2024 <1 %  
Publication
- 
- 25 Hamed Khosravi, Babak Amiri, Navid  
Yazdanjue, Vahide Babaiyan. "An improved  
group teaching optimization algorithm based  
on local search and chaotic map for feature  
selection in high-dimensional data", Expert  
Systems with Applications, 2022 <1 %  
Publication
- 
- 26 Heba Askr, Mahmoud Abdel-Salam, Aboul Ella  
Hassanien. "Copula entropy-based golden  
jackal optimization algorithm for high-  
dimensional feature selection problems",  
Expert Systems with Applications, 2024 <1 %  
Publication
-

- 27 Peda Baliyarasimhula Aakash Sai Raj,  
Jayashree Piri, Sai Babu Eluri, Sankeerth  
Reddy S. "Work Visa Analysis using Machine  
Learning Techniques", 2023 Third  
International Conference on Artificial  
Intelligence and Smart Energy (ICAIS), 2023  
Publication <1 %
- 
- 28 [ouci.dntb.gov.ua](http://ouci.dntb.gov.ua) <1 %  
Internet Source
- 
- 29 [vdoc.pub](http://vdoc.pub) <1 %  
Internet Source
- 
- 30 "Computational Intelligence in Data Mining",  
Springer Science and Business Media LLC,  
2017 <1 %  
Publication
- 
- 31 Alok Kumar Shukla, Pradeep Singh, Manu  
Vardhan. "Gene selection for cancer types  
classification using novel hybrid  
metaheuristics approach", Swarm and  
Evolutionary Computation, 2020  
Publication <1 %
- 
- 32 Anupama Namburu, Soubhagya Sankar  
Barpanda. "Recent Advances in Computer  
Based Systems, Processes and Applications",  
CRC Press, 2020 <1 %  
Publication
-

- 33 Arman Arefi, Barbara Sturm, Gardis von Gersdorff, Abozar Nasirahmadi, Oliver Hensel. "Vis-NIR hyperspectral imaging along with Gaussian process regression to monitor quality attributes of apple slices during drying", LWT, 2021 **<1 %**  
Publication
- 
- 34 Hancheng Huang, Qingwei Liang, Shanshan Hu, Cheng Yang. "Chaotic heuristic assisted method for the search path planning of the multi-BBUG cooperative system", Expert Systems with Applications, 2024 **<1 %**  
Publication
- 
- 35 Mingwei Wang, Chunming Wu, Lizhe Wang, Daxiang Xiang, Xiaohui Huang. "A feature selection approach for hyperspectral image based on modified ant lion optimizer", Knowledge-Based Systems, 2019 **<1 %**  
Publication
- 
- 36 Mohammed A. Awadallah, Mohammed Azmi Al-Betar, Malik Shehadeh Braik, Abdelaziz I. Hammouri, Iyad Abu Doush, Raed Abu Zitar. "An enhanced binary Rat Swarm Optimizer based on local-best concepts of PSO and collaborative crossover operators for feature selection", Computers in Biology and Medicine, 2022 **<1 %**  
Publication
-

- 37 Nitesh Singh Malan, Shiru Sharma. "Feature selection using regularized neighbourhood component analysis to enhance the classification performance of motor imagery signals", Computers in Biology and Medicine, 2019 <1 %  
Publication
- 
- 38 Zahra Beheshti. "BMPA-TVSinV: A Binary Marine Predators Algorithm using time-varying sinus and V-shaped transfer functions for wrapper-based feature selection", Knowledge-Based Systems, 2022 <1 %  
Publication
- 
- 39 dokumen.pub <1 %  
Internet Source
- 
- 40 evo-ml.com <1 %  
Internet Source
- 
- 41 file.techscience.com <1 %  
Internet Source
- 
- 42 www.science.gov <1 %  
Internet Source
- 
- 43 Gurdeep Singh, Urvinder Singh. "Hybrid binary grey wolf naked mole-rat algorithm for fragment-type UWB antenna optimization using time-varying transfer functions", Expert Systems with Applications, 2023 <1 %  
Publication

- 44 Mohadeseh Ghayekhloo, Ahmad Nickabadi. "CLP-GCN: Confidence and label propagation applied to Graph Convolutional Networks", *Applied Soft Computing*, 2023  
Publication <1 %
- 45 Mohammadali Ahmadi. "Machine learning", Elsevier BV, 2024  
Publication <1 %
- 46 Priyanka Anand, Sankalap Arora. "A novel chaotic selfish herd optimizer for global optimization and feature selection", *Artificial Intelligence Review*, 2019  
Publication <1 %
- 47 "Artificial Intelligence and Data Science", Springer Science and Business Media LLC, 2022  
Publication <1 %
- 48 Ke Chen, Feng-Yu Zhou, Xian-Feng Yuan. "Hybrid particle swarm optimization with spiral-shaped mechanism for feature selection", *Expert Systems with Applications*, 2019  
Publication <1 %
- 49 Tansel Dokeroglu, Ayça Deniz, Hakan Ezgi Kiziloz. "A robust multiobjective Harris' Hawks Optimization algorithm for the binary classification problem", *Knowledge-Based Systems*, 2021  
Publication <1 %

## CERTIFICATE-1

Sixth International Conference on  
**Advances in Electrical and Computer  
Technologies 2024 (ICAECT 2024)**

26 - 27, September 2024 | Tiruchengode, India | [www.icact.co.in](http://www.icact.co.in)

## CERTIFICATE

SPCS 3067

Peer Reviewed

Publication Partner



This certificate is presented to



**Lakshmi Prasanna Bommu**

Department of Computer Science and Engineering,  
Narasaraopeta Engineering College, Narasaraopet,  
Andhra Pradesh, India

for presenting the research paper entitled “ Classification and Feature Selection Method for Medical Datasets by BGEO TVFL (Binary golden eagle optimization-Time Varying Flight length) and KNN (k-nearest neighbour) ” in the Sixth International Conference on Advances in Electrical and Computer Technologies 2024 (ICAECT 2024). ICAECT 2024 is jointly organized by the Senguntha Engineering College (Autonomous), Tiruchengode, TamilNadu, India and Diligentec Solutions, Coimbatore, Tamil Nadu, India during 26 – 27. September 2024. The Conference has been organized in ONLINE MODE.

Industry Partner

**DILIGENTEC SOLUTIONS**

The handwritten signature of Dr. Thangaprakash Sengodan, Conference Chair.

Dr. Thangaprakash Sengodan

Conference Chair



ICAECT

## CERTIFICATE-2

Sixth International Conference on  
**Advances in Electrical and Computer  
Technologies 2024 (ICAECT 2024)**

26 - 27, September 2024 | Tiruchengode, India | [www.icact.co.in](http://www.icact.co.in)

## CERTIFICATE

**SPCS 3067**

Peer Reviewed

This certificate is presented to



**Swetha Lakshmi Paruchuri**

Department of Computer Science and Engineering,  
Narasaraopeta Engineering College, Narasaraopet,  
Andhra Pradesh, India

Publication Partner



for presenting the research paper entitled " Classification and Feature Selection Method for Medical Datasets by BGEO TVFL (Binary golden eagle optimization-Time Varying Flight length) and KNN (k-nearest neighbour) " in the Sixth International Conference on Advances in Electrical and Computer Technologies 2024 (ICAECT 2024). ICAECT 2024 is jointly organized by the Sengunthar Engineering College (Autonomous), Tiruchengode, TamilNadu, India and Diligentec Solutions, Coimbatore, Tamil Nadu, India during 26 – 27. September 2024. The Conference has been organized in ONLINE MODE.

Industry Partner  
**DILIGENTEC SOLUTIONS**

  
Dr. Thangaprakash Sengodan  
Conference Chair



## CERTIFICATE-3

Sixth International Conference on  
**Advances in Electrical and Computer  
Technologies 2024 (ICAECT 2024)**

26 - 27, September 2024 | Tiruchengode, India | [www.icact.co.in](http://www.icact.co.in)

## CERTIFICATE

**SPCS 3067**

Peer Reviewed

Publication Partner



This certificate is presented to



**Sahithi Mellacheruvu**

Department of Computer Science and Engineering,  
Narasaraopeta Engineering College, Narasaraopet,  
Andhra Pradesh, India

for presenting the research paper entitled “ Classification and Feature Selection Method for Medical Datasets by BGEO TVFL (Binary golden eagle optimization-Time Varying Flight length) and KNN (k-nearest neighbour) ” in the Sixth International Conference on Advances in Electrical and Computer Technologies 2024 (ICAECT 2024). ICAECT 2024 is jointly organized by the Sengunthar Engineering College (Autonomous), Tiruchengode, TamilNadu, India and Diligentec Solutions, Coimbatore, Tamil Nadu, India during 26 – 27, September 2024. The Conference has been organized in ONLINE MODE.

Industry Partner

**DILIGENTEC SOLUTIONS**

A handwritten signature in black ink.

Dr. Thangaprakash Sengodan  
Conference Chair



ICAECT

## CERTIFICATE-4

Sixth International Conference on  
**Advances in Electrical and Computer  
Technologies 2024 (ICAECT 2024)**

26 - 27, September 2024 | Tiruchengode, India | [www.icaect.co.in](http://www.icaect.co.in)

## CERTIFICATE

**SPCS 3067**

Peer Reviewed

Publication Partner



This certificate is presented to



**Rama Krishna Eluri**

Department of Computer Science and Engineering,  
Narasaraopeta Engineering College, Narasaraopet,  
Andhra Pradesh, India

for presenting the research paper entitled “ Classification and Feature Selection Method for Medical Datasets by BGEO TVFL (Binary golden eagle optimization-Time Varying Flight length) and KNN (k-nearest neighbour) ” in the Sixth International Conference on Advances in Electrical and Computer Technologies 2024 (ICAECT 2024). ICAECT 2024 is jointly organized by the Sengunthar Engineering College (Autonomous), Tiruchengode, TamilNadu, India and Diligentec Solutions, Coimbatore, Tamil Nadu, India during 26 – 27, September 2024. The Conference has been organized in ONLINE MODE.

Industry Partner

**DILIGENTEC SOLUTIONS**

Dr. Thangaprakash Sengodan  
Conference Chair



## UDAYAMOSTAV CERTIFICATE

