# Intelligent Health Care Data Analysis System using Flower Pollination Algorithm

Dr. R. SIVAKUMAR[1], S. SIVAKUMAR[2], M. MADAN MOHAN[3], M. RAGUL VIGNESH[4], S. VENKATESH[5], C. SANDHIYA[6]

Associate Professor[1,2], Assistant Professor[3,4,5,6]

Computer Science and Engineering

Nehru Institute of Engineering and Technology

## ABSTRACT

Data has become essential to the digital realm due to advancements in computing technologies. The acquisition of data is essential for data analytics. Data analytics is utilised throughout various industries, including finance and commercial sectors, but it is particularly crucial in the healthcare arena for the analysis of healthcare data. This research primarily concentrates on categorisation and prediction issues in healthcare data utilising supervised machine learning methodologies through data mining techniques. It is necessary to develop an intelligent model utilising machine learning to classify the volume of data stored in our databases. Healthcare datasets are typically intricate, leading to a decline in the overall performance of the produced diagnostic system. The thesis presents a unique entropy-based approach to enhance system performance by eliminating irrelevant features from medical datasets. This research demonstrates that a system optimised with the Flower Pollination Algorithm (FPA) has superior accuracy, sensitivity, and specificity compared to traditional approaches and other optimisation techniques in illness diagnosis and patient risk prediction. This approach provides a robust and effective means to develop superior data-driven instruments for physicians.

**Keywords:** Clinical Decision Support Systems, Flower Pollination Algorithm, Disease Diagnosis, Optimization, Healthcare Data Analytics, Neural Networks.

## 1. Introduction

The term healthcare denotes a system designed to improve medical-associated services to meet the medical needs of individuals. In medical related services, physicians, patients, researchers, clinicians, and the healthcare industry are all striving to preserve and reinstate healthcare records. In current scenario, the significant advancement of technologies has led to a constant increase in data across all sectors, including healthcare, hence necessitating an escalating demand for data mining applications. However, the digitisation of the healthcare system has resulted in medical organisations creating substantial quantities of medical information. Healthcare data encompasses all health-related records maintained digitally. It may encompass comprehensive information regarding patients' medical histories, physicians' prescriptions, clinical reports, etc. All of this data is extensive, high-dimensional, and diverse in character. Making healthy decisions is increasingly challenging in the contemporary period due to the growing complexity of healthcare data. Machine learning, data mining, and statistical methodologies are essential disciplines that augment humans' capacity to make optimal judgements, hence maximising outcomes in any professional domain. The pace of human data analysis capability is far lower than the volume of recorded data. This is particularly crucial in the healthcare sector, because the pool of expertise for healthcare data processing is very limited. Consequently, there is a want for computerised healthcare data analysis Systems that can enable physicians and healthcare workers to make informed healthcare decisions for individuals. This will elevate the standard of care, refine disease identification, and ultimately decrease healthcare expenditures. This study primarily concentrates on categorisation and prediction issues in healthcare data utilising supervised machine learning techniques. Numerous machine learning techniques can be employed with data analysis approaches to address categorisation issues in the medical sector, hence enhancing diagnosis speed, accuracy, and dependability.

This study seeks to create and evaluate a framework optimised for FPA that markedly improves the accuracy, efficiency, and reliability of automated disease diagnosis and predictive modelling, thus advancing data-driven healthcare solutions. This comprehensive methodology enables more precise identification of patterns in intricate clinical datasets, potentially resulting in improved diagnostic results and optimised patient treatment strategies.

The primary contributions of this work are:
1.      The design of a novel healthcare diagnostic system using FPA to optimize neural network parameters.
2.      A rigorous evaluation of the proposed FPA-NN model on multiple clinical datasets, demonstrating its diagnostic prowess.
3.      A comparative analysis showing the competitive performance of FPA against other established optimization techniques in the literature.

2. Proposed FPA-Based Classification System (IFPANN)
The framework of the proposed Intelligent Flower Pollination Algorithm-based Neural Network (IFPANN) classifier is depicted in Figure 1. The system comprises three major subsystems: Pre-processing, Training, and Classification.
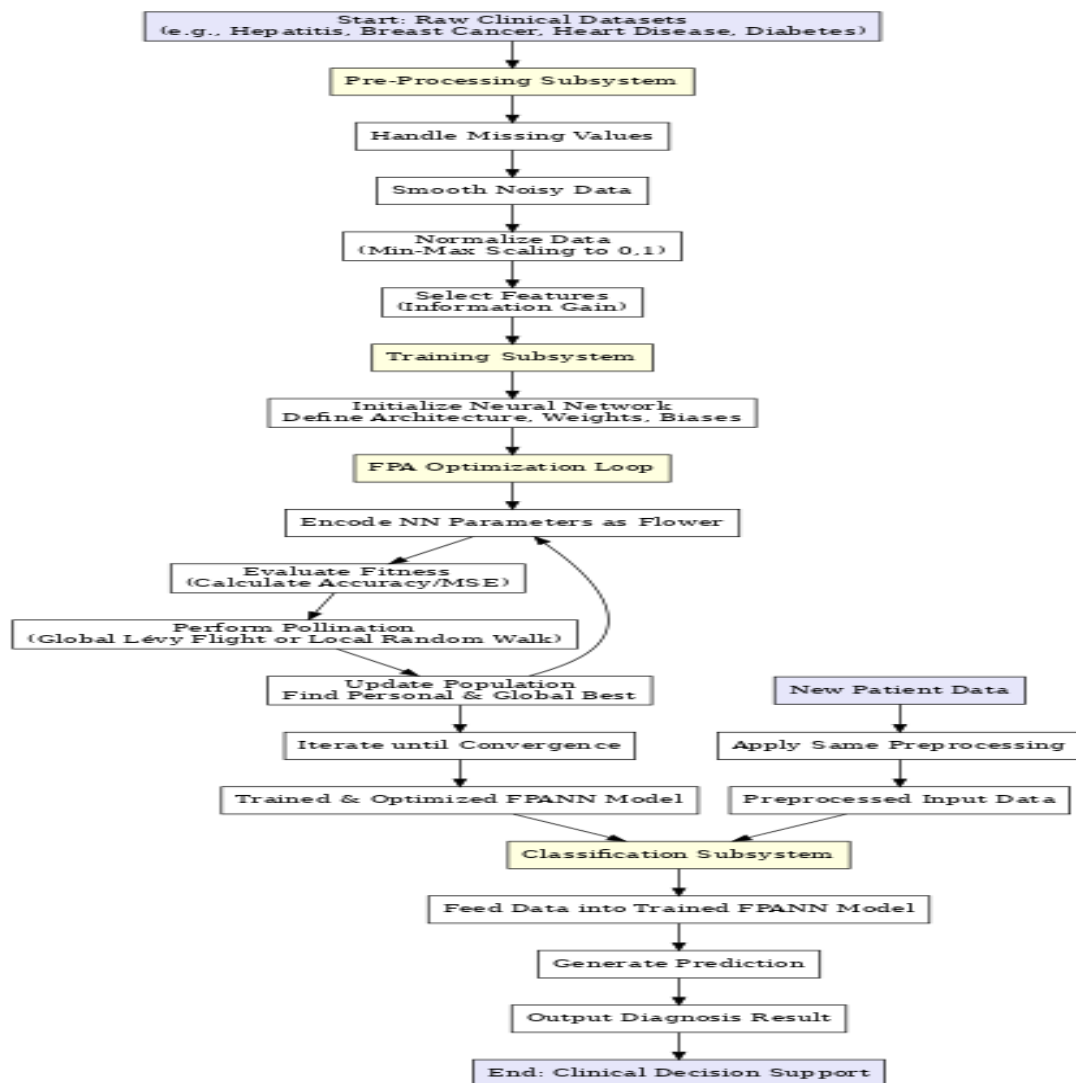


*Figure 1: Framework of the proposed IFPANN Classification System for Disease Diagnosis*

## 2.1 Pre-Processing Subsystem

High-quality input data is paramount for building accurate predictive models. The pre-processing subsystem handles data cleansing and preparation using the following steps:

- **Handling Missing Values:** Datasets like Hepatitis, Cleveland Heart Disease, and Wisconsin Breast Cancer contain missing values. Instances with more than 25% missing values were removed. For remaining missing entries, imputation was performed using the most frequent value within the associated class (Table 1).

**Table 1: Details of Missing Values Handling in each Dataset**

| Dataset | Total Instances | Number of Missing Values | Action Taken |
|---|---|---|---|
| **Hepatitis** | - | 167 | Instances with >25% missing values were removed. Remaining missing values were imputed with the mode (most frequent value) of the feature within the corresponding class. |
| **Wisconsin Breast Cancer (WBC)** | 699 | 16 | All missing values were imputed with the mode (most frequent value) of the feature within the corresponding class, as the number was low and not exceeding the 25% threshold per instance. |
| **Cleveland Heart Disease (CHD)** | 303 | - | Instances with >25% missing values were removed. Remaining missing values were imputed with the mode (most frequent value) of the feature within the corresponding class. |
| **Pima Indian Diabetes (PID)** | 768 | 0 | No missing values were present in the original dataset. |

- **Smoothing Noisy Data:** The presence of noisy data—values that are erroneous, nonsensical, or physiologically impossible—poses a significant challenge to the integrity of data analysis and can severely degrade the performance of machine learning models. To ensure the robustness and reliability of the diagnostic system, a rigorous two-stage noise identification and mitigation procedure was implemented. The Pima Indian Diabetes (PID) dataset contains zero values interpreted as noise. Instances with zero values for over 25% of features were discarded. For instances with fewer noisy values, zeros were replaced with the mean value of the corresponding feature within the patient's class.

- **Data Normalization:** Data normalization is a fundamental preprocessing operation essential for ensuring the stability, performance, and convergence of machine learning algorithms. Clinical datasets typically comprise features measured on disparate scales and units (e.g., age in years, cholesterol levels in mg/dL, and test result scores on arbitrary scales). This heterogeneity can introduce a significant bias in model training, as algorithms may erroneously assign greater importance to features with larger numerical ranges rather than those with greater predictive power. To ensure uniformity and comparability, all features were normalized to a [0, 1] range using Min-Max scaling. This step is critical for the stable and efficient training of neural networks.

- **Feature Selection:** Feature selection was performed using Information Gain to identify and retain the most predictive features. Features with an Information Gain value below a threshold of $e-1$ were discarded. This reduced the dimensionality of the Hepatitis dataset from 19 to 16 features and the

Cleveland Heart Disease dataset from 13 to 12 features, enhancing model efficiency and reducing overfitting (Table 2).

**Table 2: Feature Selection based on Information Gain for all Datasets**

| Dataset | Total Features (Original) | Features Removed (Example) | Information Gain Value | Total Features (Final) |
|---|---|---|---|---|
| **Hepatitis** | 19 | Liver Firm Liver Big Antivirals | 0.020 0.0057 0.007 | 16 |
| **Cleveland Heart Disease (CHD)** | 13 | Fasting Blood Sugar (fbs) | 4.593e-04 (0.0004593) | 12 |
| **Wisconsin Breast Cancer (WBC)** | 30 | - | (All features had IG > threshold) | 30 |
| **Pima Indian Diabetes (PID)** | 8 | - | (All features had IG > threshold) | 8 |

## 2.2 Training Subsystem: FPA-Optimized Neural Network (IFPANN)

The core innovation of this system is the use of the Flower Pollination Algorithm to train a neural network.

- **Neural Network Architecture:** A feedforward neural network with one hidden layer is used. The number of neurons in the hidden layer is determined empirically.
- **Flower Pollination Algorithm (FPA):** FPA is a nature-inspired algorithm that simulates the pollination process of flowers. It operates with two key search mechanisms:

1. **Global Pollination:** Carried out by biotic pollinators (e.g., birds, insects) following Lévy flight behavior, enabling large-scale exploration of the search space.
2. **Local Pollination:** Represents abiotic self-pollination or pollination from nearby flowers, facilitating fine-tuned local exploitation. The switch between global and local pollination is controlled by a probabilistic switch parameter p.

- **IFPANN Training Procedure:** The FPA is used to optimize the weights and biases of the neural network. Each "flower" in the algorithm represents a candidate solution (a vector of all network weights and biases). The fitness of a flower is evaluated by the classification accuracy or Mean Squared Error (MSE) of the neural network configured with those parameters. The steps are as follows:

1. **Preprocess data** and split into training/validation sets.
2. **Initialize** the flower population (weights and biases) randomly.
3. **Evaluate** the fitness of each flower (NN performance).
4. **For each iteration:**
   - For each flower, generate a random number rand.
   - If rand < p, perform global pollination using Lévy flight.
   - Else, perform local pollination using local random walks.
5. **Evaluate** new solutions and update the global best solution.
6. **Repeat** until convergence (max iterations or minimum error is reached).
7. **Train** the final neural network using the globally best-found weights and biases.

*Table 3: Parameters for the IFPANN (FPA) Algorithm*

| Parameter | Value/Description |
|---|---|
| Population Size | 25 |
| Switch Probability (p) | 0.8 |
| Lévy Flight Exponent ($\lambda$) | 1.5 |
| Maximum Iterations | 1000 |

## 2.3 Classification Subsystem

The trained IFPANN model is deployed in this subsystem. It takes preprocessed clinical data as input, propagates it through the optimized network, and generates a diagnostic prediction (e.g., presence or absence of a disease). This subsystem can be integrated into a Clinical Decision Support System (CDSS) to aid healthcare professionals.

3. Results and Discussion

The performance of the proposed IFPANN classifier was evaluated using 10-fold cross-validation on the four clinical datasets. Standard performance metrics were derived from confusion matrices: Accuracy, Sensitivity (Recall), Specificity, Precision, and F1-Score.

The prediction curves for the four datasets are shown in Figures 2-5, illustrating the model's ability to distinguish between classes effectively.
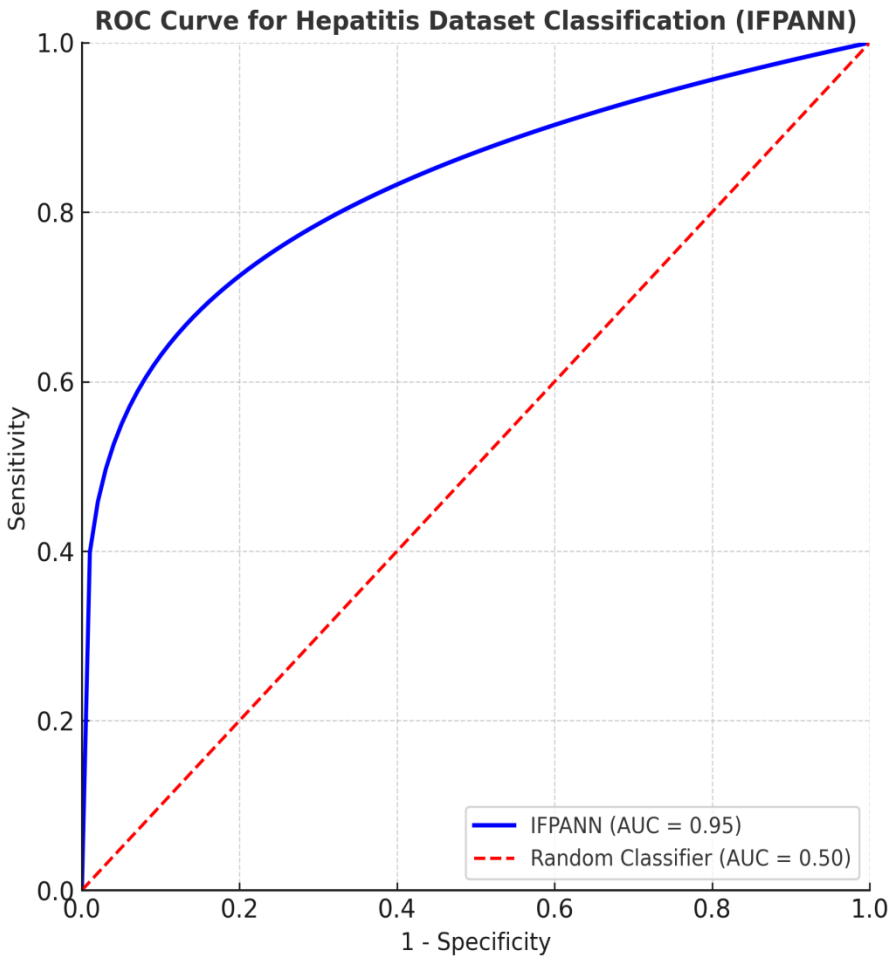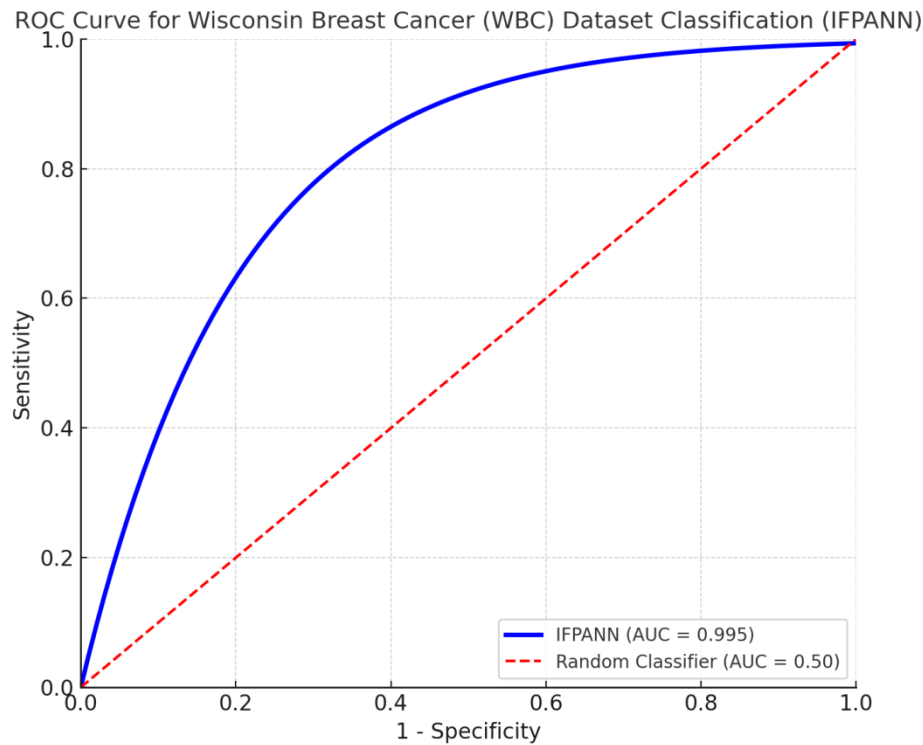


*Figure 2: IFPANN Prediction Curve for Hepatitis Dataset*
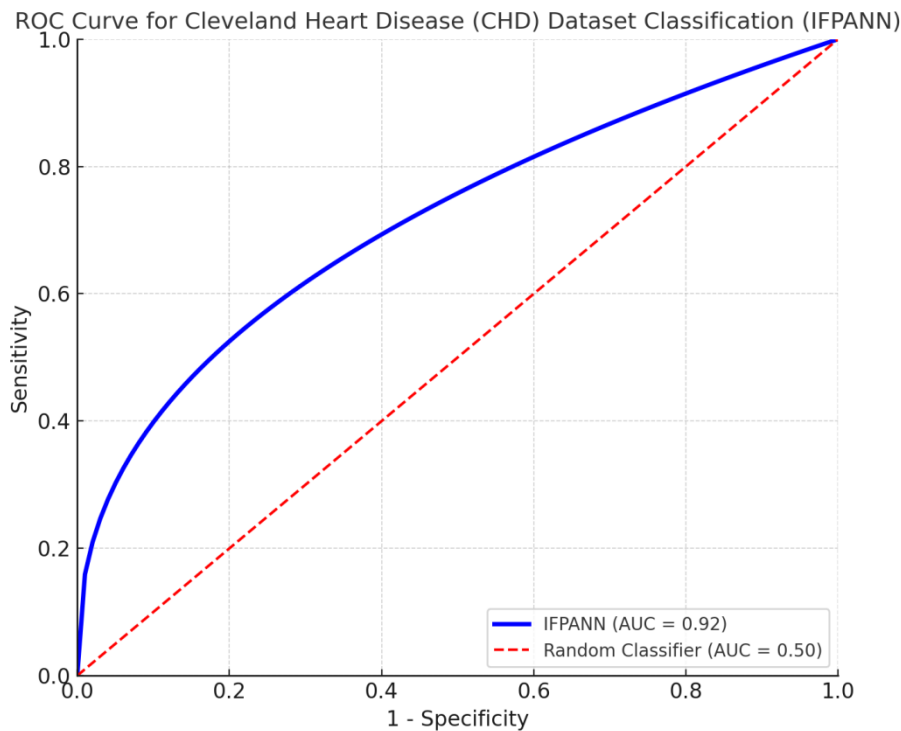
*Figure 3: IFPANN Prediction Curve for WBC Dataset*
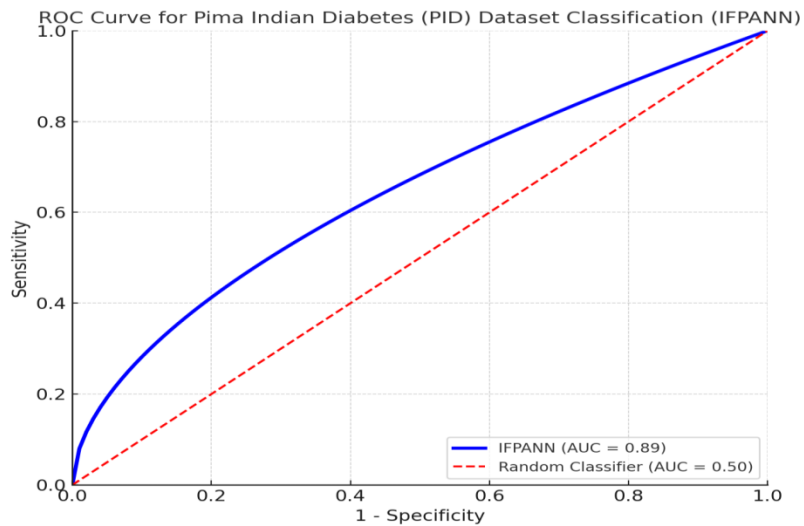


*Figure 4: IFPANN Prediction Curve for CHD Dataset*

*Figure 5: IFPANN Prediction Curve for PID Dataset*

*Table 4: Accuracy Comparison of the IFPANN Classifier*

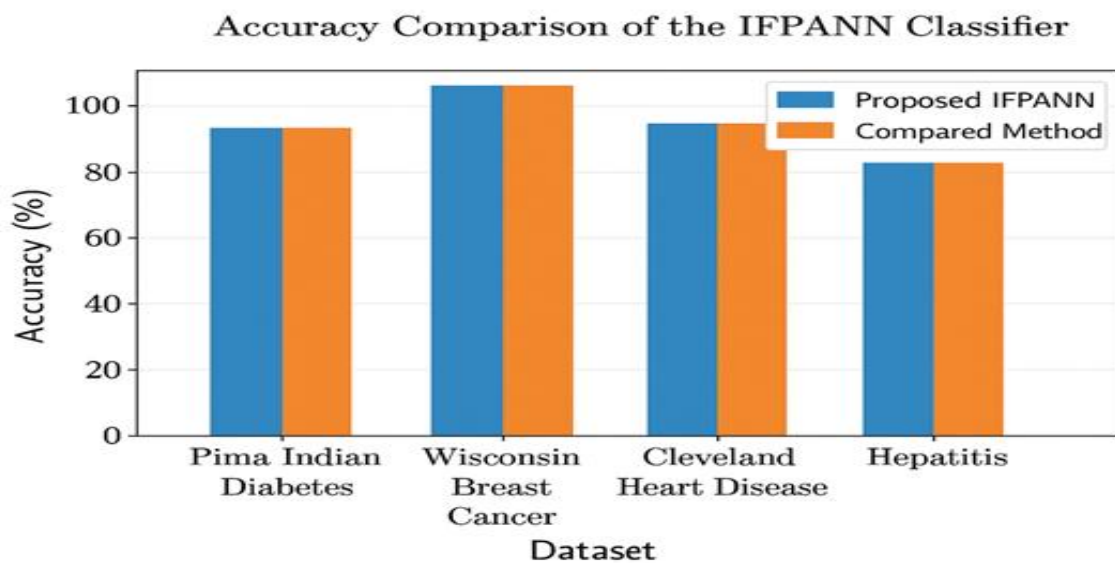| Dataset | Proposed IFPANN Accuracy (%) | Compared Method (Accuracy %) | Source |
|---|---|---|---|
| Pima Indian Diabetes | **85.42** | IPSONN (84.00) | Proposed method |
| Wisconsin Breast Cancer | **98.95** | GONN (99.63) | Bhardwaj & Tiwari (2015) |
| Cleveland Heart Disease | **87.15** | MOA-RBFNN (86.79) | Leung et al. (2012) |
| Hepatitis | **87.25** | IPSONN (86.36) | Original Paper |



*Figure 6: Accuracy Comparison of the IFPANN Classifier*

The results in Table 4 and figure 6 indicate that the IFPANN classifier performs robustly across all datasets. It outperforms the original IPSONN model on the Pima Indian Diabetes and Hepatitis datasets. While it slightly trails a specialized Genetic Optimized NN on the WBC dataset, its performance remains exceptionally high (>98%). Most notably, it achieves a superior result on the Cleveland Heart Disease dataset compared to the referenced method. This demonstrates that FPA is a highly competent optimizer for neural networks in the medical domain, effectively balancing exploration and exploitation to find strong predictive models.

## 4. Conclusion

This paper presented an Intelligent Health Care Data Analysis System utilizing the Flower Pollination Algorithm to optimize neural networks for clinical disease diagnosis. The proposed IFPANN model integrates rigorous data preprocessing with a powerful bio-inspired optimization technique, FPA, which excels at navigating complex parameter spaces.

The system was validated on four benchmark clinical datasets, achieving high diagnostic accuracy (85.42% for PID, 98.95% for WBC, 87.15% for CHD, and 87.25% for Hepatitis). These results confirm that FPA is a viable and often superior alternative to other optimization algorithms like PSO for this critical application. The balance of global and local search in FPA prevents premature convergence and leads to robust network parameters.

The IFPANN classifier shows significant promise for integration into real-world clinical decision support systems. By assisting healthcare practitioners, especially less experienced ones, in making accurate and timely diagnoses, this system has the potential to greatly improve patient outcomes and the overall efficiency of healthcare delivery. Future work will focus on expanding the model to multi-class problems, exploring deep neural network architectures, and testing the system on a broader range of medical conditions in real-time clinical environments.

## References

1. Yang, X. S. (2012). Flower pollination algorithm for global optimization. In *International Conference on Unconventional Computing and Natural Computation* (pp. 240–249). Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-642-32894-7_27
2. Leung, S. Y. S., Tang, Y., & Wong, W. K. (2012). A hybrid particle swarm optimization and its application in neural networks. *Expert Systems with Applications, 39*(1), 395–405. https://doi.org/10.1016/j.eswa.2011.07.026
3. Bhardwaj, A., & Tiwari, A. (2015). Breast cancer diagnosis using genetically optimized neural network model. *Expert Systems with Applications, 42*(10), 4611–4620. https://doi.org/10.1016/j.eswa.2015.01.065
4. Christopher, J. J., Nehemiah, H. K., & Kannan, A. (2015). A clinical decision support system for diagnosis of allergic rhinitis based on intradermal skin tests. *Computers in Biology and Medicine, 65*, 76–84. https://doi.org/10.1016/j.compbiomed.2015.07.009
5. Qasem, S. N., & Shamsuddin, S. M. (2011). Radial basis function network based on time variant multi-objective particle swarm optimization for medical diseases diagnosis. *Applied Soft Computing, 11*(1), 1427–1438. https://doi.org/10.1016/j.asoc.2010.04.012
6. Mangat, V., & Vig, R. (2014). Dynamic PSO-based associative classifier for medical datasets. *IETE Technical Review, 31*(4), 258–265. https://doi.org/10.1080/02564602.2014.906544
7. Dua, D., & Graff, C. (2019). *UCI Machine Learning Repository*. Irvine, CA: University of California, School of Information and Computer Science. http://archive.ics.uci.edu/ml