# A New hybrid Genetic Search Algorithm and Invasive Weed Optimization Algorithms for Skin Lesion Cancer Classification

Dr. S. M. Uma[1], Dr. D. Sivakumar[2]

[1]Associate Professor,Department of CSE, Kings College of Engineering, Thanjavur, Tamil Nadu, India

[2]Assistant Professor, Department of CSE, Kings College of Engineering, Thanjavur, Tamil Nadu, India

*Abstract* - **Skin disease is a primary hassle amongst people global. Different learning algorithm getting to know. Strategies can be implemented to perceive lessons of pores and skin sickness. Accurately diagnosing skin lesions to discriminate among benign and malignant skin lesions is critical to make certain suitable affected person treatment. Skin malignant growth is one of most dangerous maladies in people. As per the high closeness among melanoma and nevus sores, doctors set aside substantially more effort to explore these sores. This paper displays another technique dependent on enhancement calculation to order and foresee skin malignant growth maladies tried utilizing certifiable disease datasets. This philosophy going to joins new two sort of calculation. One is genetic algorithm(GA) and another is Invasive weed optimization (IWO)algorithm to arrange and anticipate malignant growth prior. The proposed framework is assessed by arranging and expectation malignant growth sicknesses in skin sore disease datasets and assessment measures. The outcomes are thought about with (convolution algorithm)SVM execution benchmark. Framework can defeat to diagnosing the malady rapidly and exactness. Contrasting with other calculation proposed calculation has more precision.**

*Key Words: IWO, SVM, SSA data set, Analysis, Clustering, Accuracy*

## 1. INTRODUCTION

Information mining is the procedure where esteemed data is separated from the enormous dataset. It has arrived at the high development over recent years. Because of the convenience of information mining approaches in wellbeing world, it has become the great innovation in medicinal services area. Malignant growth is a speculatively last ailment caused fundamentally by conservational issues that change qualities encoding basic cell administrative proteins. Resultant Many highlights of the cutting edge Western eating routine (high fat, low fiber content) will expand malignant growth recurrence.

## 2. METHODOLOGY SYSTEM IMPLEMENTATION

The following actions are carried out in the proposed system. They are;

1. Dataset Acquisition
2. Preprocessing
3. Feature Selection
4. Disease Diagnosis
5. Evaluation Criteria

### DATASET ACQUISITION

In this module, transfer the datasets. The dataset might be microarray dataset. Accumulate the information from emergency clinics, server farms and disease inquires about focuses. The gathered information is pre-handled and put away in the information base to fabricate the model.

### PREPROCESSING:

Data pre-handling is a significant advance in the information mining process. The expression "manure in, trash out" is for the most part relevant to information mining and machine ventures. Information gathering strategies are regularly shakily controlled, coming about in out-of-go values, inconceivable information blend, missing qualities, and so forth. Investigating information that has not been deliberately screened for such issues can deliver equivocal outcomes.

### FEATURE SELECTION:

In this module is utilized to choose the highlights of the given dataset. Credit choice was performed to decide the subset of highlights that were exceptionally related with the class while having low inter correlation.

### DISEASE DIAGNOSIS:

Based on the values acquired from training phase, the performance of the NN network is analyzed to obtain appropriate values for testing phase. In order to find the optimum structure, the NN network performance has been analyzed for the optimum number of hidden nodes and epochs. For this situation, the epochs will be set to a definite preset value. Then, the NN network was trained at the appropriate range of hidden nodes. The number of hidden

**Special Issue - 2021**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**ICRADL - 2021 Conference Proceedings**

nodes that have given the best performance is then selected as the optimum hidden nodes. After that, by fixing the optimum number of hidden nodes, the epochs will be analyzed in a similar way to obtain the optimal number of periods that container give the highest or best accuracy.

EVALUATION CRITERIA:

In this module, the exhibition of the proposed Genetic calculation is broadly examined with that of some current directed and solo quality bunching and quality determination calculations utilizing different factual measures.

DATA COLLECTION:

We have collected our dataset of signs and symptoms and patient details related to ten not unusual skin sicknesses from the department of Skin and any actual time Hospital and dialogues with concerned health practitioner.

2.4.1 INPUT DATASET (NUMERIC):

The input dataset for this project are contain attribute of various skin diseases. These attribute were taken as input from the data set.

PRE-PROCESSING:

In pre-processing practice Gaussian filtering to our input photograph. Gaussian filtering is often used to get rid of the noise from the image. Here we used wiener function to our input picture. Gaussian clear out is windowed filter out of linear class, by using its nature is weighted imply. Named after well-known scientist Carl Gauss due to the fact weights within the clear out calculated according to Gaussian distribution. feature.

CLASSIFICATION :

Classification is a records mining feature that assigns gadgets in a group to target categories or lessons. The aim of category is to accurately are expecting the target class for each case in the data. Classification divides records samples into target classes. The class technique predicts the goal elegance for each records points.

### 3. SUPPORT VECTOR MACHINE (SVM):

SVM is a gathering of learning calculations essentially utilized for arrangement undertakings on confounded information, for example, picture characterization and protein structure examination. SVM is utilized in innumerable fields in science and industry, including Bio-innovation, Medicine, Chemistry and Computer Science. It has additionally ended up being preferably appropriate for order of huge content archives, for example, those housed in for all intents and purposes all huge, present day associations. Presented in 1992, SVM rapidly moved toward becoming viewed as the best in class strategy for order of mind boggling, high-dimensional information. Specifically its capacity to catch patterns seen in a little preparing set and to

sum up those patterns against a more extensive corpus have made it helpful over countless. SVM utilizes a managed learning approach, which implies it figures out how to group inconspicuous information in view of a lot of namedpreparing information, for example, corporate records. The underlying arrangement of preparing information is ordinarily recognized by space specialists and is utilized to construct a model that can be connected to some other information outside the preparation set. The exertion required to develop a great preparing set is very unassuming, especially when contrasted with the volume of information that might be at last ordered against it. This implies learning calculations, for example, SVM offer an outstandingly practical technique for content characterization for the enormous volumes of archives created by present day associations. The equalization of this paper covers the inward functions of SVM, its application in science and industry, the lawful faultlessness of the technique just as order precision contrasted with manual characterization.

The Classifier are currently connected utilizing the highlights processed in the past sub-area that incorporates coiflet wavelet change of the EMG flags and highlights like mean, vitality and standard deviation for the D4 coefficient of the change. These highlights are currently used to prepare the SVM model and after that test on the equivalent prepared SVM for the characterization. SVM (Support Vector Machine) is a classifier which requires preparing of the model for the testing of the examples and such procedure in which preparing is given is called as directed learning strategy. In this, demonstrate is right off the bat prepared to pick up as indicated by the highlights of the classes which should be ordered. At that point that prepared model is misused in the grouping procedure. Learning given to the classifier gives better outcomes contrasted with the different classifiers in which unsupervised learning is utilized. This classifier is straightforward and simple to comprehend as it builds a hyperplane between the distinctive classes which should be ordered. The classifier utilized might be straight or non-direct. In straight SVM, preparing tests of the classes are straightly distinguishable. Yet, it is troublesome in down to earth circumstances that a straight line is adequate to arrange every single example. For such cases, non-straight classifier is abused.

### 4. ESSENTIAL ALGORITHMS

#### 4.1 THE PROPOSED GA APPRAOCH
This section describes a GA that evolves a population of individuals, where each individual represents a classification rule. More precisely, each individual represents the antecedent (If part) of a classification rule. The consequent (Then part) of the rule is not encoded in the genome. Rather, it is fixed for a given GA run, so that in each run all the individuals represent rules with the same consequent. The remaining details of GA are as given below.

Individual Representation
In the proposed algorithm the genome of an individual consists of a conjunction of conditions composing a given rule antecedent. An individual is encoded as a set of $n$

conditions, where *n* is the number of predictor attributes(Figure 1). Each gene represents a rule condition of the form $A_i \, Op_i \, V_{ij}$, where[4]:
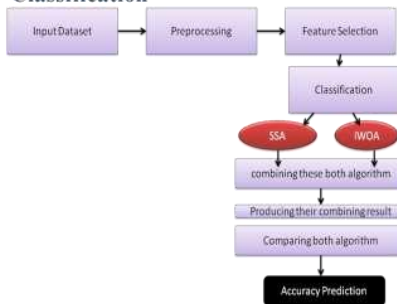
- $A_i$ denotes the *i*-th predictor attribute;

- $Op_i$ denotes the comparison operators-e.g., "=" for categorical attributes; ",<=" or ">" for continuous attributes-used in the *i*-th conditions;

- $Val_{ij}$ denotes the *j*-th value of the domain of $Attr_i$ .

In this example the data being mined has only three attributes, so three conditions are encoded in the genotype . This genotype is decoded into the following rule antecedent If *Marital-status= married ☐ Age > 30*. The consequent (Then part) of the rule, which specifies the predicted class, is not represented in the genome. Rather, it is fixed for a given GA run, so that all individuals have the same rule consequent during all that run.

### 4.2 Genetic Operators
Genetic operators are being used in GAs to maintain genetic diversity by introducing new genetic material and to manipulate or recombine the genetic material of candidate rules. In the proposed system tournament selection, with tournament size of 5, two-point crossover with crossover probability of 90%, and mutation probability of 1% are used.



Furthermore, we use elitism with an elitist factor of 1 - i.e. the best individual of each generation is passed unaltered into the next generation. The crossover (or recombination) operator essentially swaps genetic material between two "parents" creating two new "child" individuals. These two crossover points are randomly generated such that the genes between these two points are swapped between the two individuals, yielding the new child individuals. Note that crossover points can fall only between genes, and not inside a gene. Hence, crossover swaps entire rule conditions between individuals, but it cannot produce new rule conditions. The mutation is an operator that acts on a single individual at a time. It helps to find the global optimal solution of the problem, searching new areas. This operator randomly transforms the value of an attribute into another (different) value belonging to the domain of that attribute.

### 4.3 Fitness Function
The fitness function is used to evaluate how well an individual solves the target problem, and it is responsible for determining which individuals will reproduce and have parts of their genetic material passed onto the next generation. The better the fitness of an individual, the higher the probability of that individual being selected for reproduction, crossover, and mutation operations. Before we can define the fitness function, it is necessary to recall a few basic concepts on classification rule evaluation. The discovered rule(R) in the form If P Then D; where P is antecedent part and D is the consequent part.

## 5. THE PROPOSED HYBRID ALGORITHM BASED ON GA AND IWO

### 5.1 THE PROPOSED HYBRID ALGORITHM
Recently, some research groups focused attention on exploitation of GA potentialities for information extraction from biomedical database. This paper uses a hybrid GA with IWO Search (LS) for feature selection. AML-ALL dataset was used for the work.

### 5.1 DATASET

ALL and AML were based on gene expression profiling using DNA microarrays.Datasets have measurements corresponding to ALL and AML samples from Bone Marrow/peripheral blood. Leukemia (AML-ALL) included 7129 probes where 2 leukaemia variants were available: AML 25 samples; and, ALL 47 samples .

### 5.2 CLASSIFICATION OF SSA AND IWOA USING NUMERIC DATA SET
Learning where a preparation set of effectively recognized perceptions is accessible. The comparing unsupervised technique is
known as bunching, and includes gathering information into classes dependent on some proportion of comparable trademark or separation. A calculation that actualizes arrangement, particularly in a solid usage, is known as a classifier. The expression "classifier" once in a while moreover alludes to the numerical capacity.

### 5.3 CLASSIFICATION ALGORITHMS AND TUNING PARAMETERS:

Tuned parameters assume a critical job in delivering high exactness results when utilizing SVM, RF, and kNN. Each classifier has distinctive tuning steps and tuned Parameters. For each classifier, we tested a progression of qualities for the tuning procedure with the ideal parameters decided dependent on the most elevated by and large classification exactness. In this investigation, the classified results under the ideal parameters of each classifier were utilized to look at the execution of classifiers.

### 5.3.1 BOLSTE VECTOR MACHINE (SVM):

SVM classifier is generally utilized and demonstrates a decent exhibition. Along these lines, we utilized the RBF part to actualize the SVM calculation. There are two parameters that should be set while applying the SVM classifier with RBF portion: the ideal parameters of cost (C) and the part width

parameter (γ). The C parameter chooses the measure of misclassification took into consideration non-detachable preparing information, which makes the alteration of the inflexibility of preparing information conceivable. The portion width parameter (γ) influences the smoothing of the state of the class-partitioning hyperplane. Bigger estimations of C may prompt an over-fitting model, though expanding the γ esteem will influence the state of the class-partitioning hyperplane, which may influence the classification precision results. Following the investigation of Li et al. also, pretested to our dataset, in this examination, to find the ideal parameters for SVM, ten estimations of C (4−2, 2−3, 1, 21, 22, 23, 24, 15, 26, 27), and ten estimations of γ (1−5, 1−4, 1−3, 1−2, 1−1, 20, 21, 22, 23, 27) were tried. This technique was connected to each of the 14 sub-datasets.

5.3.2. RANDOM FOREST (RF):

So as to actualize the RF, two parameters should be set up the quantity of trees (ntree) and the quantity of highlights in each split (mtry). A few examinations have expressed that the agreeable outcomes could be accomplished with the default parameters. Be that as it may, as per Liaw and Wiener, the substantial number of trees will give a steady aftereffect of variable significance.

## 6. ANALYSIS OF RESULTS AND DISCUSSION:

ISSA is proposed by presenting the generation of IWO into SSA, which is contrasted and ALO, DA, PSO, IWO, what's more, SSA in the entire paper. Right now, is used to perform on 36 benchmark capacities for the base enhancement and is applied to be joined with SVM (ISSA-SVM) for the evaluation classification of air quality and with DML (ISSA-DML) for DOA estimation of MEMS vector hydrophone. The theories in SSA are that there is just a single squirrel on each tree. Right now, squirrel is respected to have its off springs. The quantity of off springs of each Eflying squirrel on the hickory nut tree, or the oak seed nut trees, or the ordinary trees of IWO is resolved, separately. At that point the number of inhabitants in the Eflying squirrels is predetermined by the method for the rising fitness values, appeared in which makes the populace assorted variety. The flying squirrels are redistributed on the hickory nut tree, or the accord nut trees, or the ordinary trees. Hence the Eflying squirrel on the hickory nut tree is kept to be the ideal in each cycle. The refreshed methodology in ISSA is performed by reobtaining the squirrel on each tree from the offspring of the squirrel on each tree, which shows the assorted variety of populace. Other than this methodology, the staying of ISSA is the equivalent as SSA. As indicated by FIGURE 5, despite the fact that the decent variety plot of IWO is diminishing with the cycles of decent variety plots of ISSA and SSA is comparative. Be that as it may, the multiplication of ISSA refreshes the populace continually, which makes ISSA have the more drawn out execution time and the bigger computational multifaceted nature. Likewise, the size of populace what's more, the quantity of cycles have an on ISSA. There are numerous enhancements for SSA to be required, such as introduction, the refreshed areas, regular checking condition and arbitrary

migration toward the finish of winter season. Moreover, more than one swarm knowledge calculation is utilized to be joined with SSA to set up the new half and half calculation.

## 7. CONCLUSION FUTURE WORK:

Two models ISSA-SVM and ISSA-DML are worked for playing out the evaluation classifications of air quality by joining ISSA with SVM and assessing the edges of DOA on reproduction analyzes by consolidating ISSA with DML, individually. By correlation, the proposed ISSA has the better combination execution and the better normal capacity esteems on 36 benchmark capacities, which shows that ISSA is able to do work advancement. At that point the proposed ISSA-SVM model has the most noteworthy classified exactness 87.91971% for figuring it out the evaluation classification of air quality, and the proposed ISSA-DML model has the least RMSE and the DOA estimations of ISSA-DML are the nearest to two sorts of episode points Thusly, the proposed calculation ISSA right now reasonable for work streamlining, and the built up models ISSA-SVM and ISSA-DML are for grade classifications of air quality and the DOA estimation, individually. These give us signs that in a future work, we will propose new or improved swarm insight calculations and apply them to advance the parameters of AI for classifications and estimation in the real world by being consolidated with different methodologies.

## REFERENCES:

[1] "Prediction of Skin Diseases using Data Mining Techniques" S. Reena Parvin1, O.A. Mohamed Jafar2.

[2] "Machine Learning Approaches to Multi-Class Human Skin Disease Detection" Ms. Seema Kolkur 1, Dr. D.R. Kalbande 2, Dr. Vidya Kharkar.

[3] Amino Singh Clair and Raminder Preet Kaur "Software Effort Estimation using k-Nearest Neighbour (kNN) "[4] Mohammad Rezwanul Huq, Ahmad Ali, Anika Rahman "Sentiment Analysis on Twitter Data using KNN and SVM"(IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 8, No. 6, 2017.

[4] Phan Thanh Noi, Martin Kappas "Comparison of Random Forest, k-Nearest Neighbor, and Support Vector Machine Classifiers for Land Cover Classification Using Sentinel-2 Imagery "

[5] Bonabeau E., Dorigo M., and Theraulaz G. (`999), "Swarm Intelligence: From natural to artificial systems", Oxford University press, NY, pp.1-25.

[6] Lillesand, T.M., Kiefer, R.W. and Chipman, J.W. (2003),"Remote Sensing and Image Interpretation", Fifth Edition,Wiley & Sons Ltd., England, pp.586-592.

[7] Ma, H. (2009), "An analysis of the behavior of migration models for Biogeography-based Optimization", Submitted for publication.

[8] Kazemian, M., Ramezani, Y., Lucas, C., Moshiri, B. (2006),"Swarm Clustering Based on Flowers Pollination by Artificial Bees", Studies in Computational Intelligence (SCI),vol. 34, Springer Berlin Publishers, New York, pp.191202.

[9] Panchal V.K., Singh, P., Kaur, N. and Kundra, H.(2009),"Biogeography based Satellite Image Classification",International Journal of Computer Science and informationSecurity, vol. 6, no.2, pp.269-274.

[10] Simon, D. (2008), "Biogeography-based Optimization", IEEE Transactions on Evolutionary Computation, vol. 12,No.6, IEEE Computer Society Press. pp. 702-713.

[11] Simon, D., Ergezer, M. and Du, D.(2009), "PopulationDistributions in Biogeography-Based OptimizationAlgorithms with Elitism", IEEE Conference on Systems,Man, and Cybernetics, San Antonio, TX, pp. 1017-1022.

**Special Issue - 2021**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**ICRADL - 2021 Conference Proceedings**

[12]   P. M. Narendra and K. Fukunaga, ``A branch and bound algorithm for feature subset selection,'' IEEE Trans. Comput., vol. C-26, no. 9, pp. 917_922, Sep. 1977.

[13]   [1]          Q. Cheng, H. Zhou, and J. Cheng, ``The Fisher-Markov selector: Fast selecting maximally separable feature subset for multiclass classication with applications to high-dimensional data,'' IEEE Trans. Pattern Anal. Mach. Intell., vol. 33, no. 6, pp. 1217_1233, Jun. 2011.

[14]   Kaur, Manpreet, Heena Gulati, and Harish Kundra. "Data mining in Agriculture on crop price prediction: Techniques and Applications." International Journal of Computer Applications 99.12 (2014): 1-3.

[15]   Wu, H., Wu, H., Zhu, M., Chen, W., & Chen, W. (2017). A new method of large-scale short-term forecasting of agricultural commodity prices: illustrated by the case of agricultural marketsin Beijing. Journal of Big Data, 4(1).

[16]   Manjula, E., & Djodiltachoumy, S. (2017). A Model forPrediction of Crop Yield. International Journal of Computational Intelligence and Informatics, 6(4), 2349 6363