# Effective and Efficient Ranking and Re-Ranking Feature Selector for Healthcare Analytics

## S.Ilangovan[1], Dr. A. Vincent Antony Kumar[2]

[1]K. L. N. College of Engineering, Madurai, India.
[2]PSNA College of Engineering and Technology, Dindigul, India

## ABSTRACT

In this work, a Novel Feature selection framework called SU embedded PSO Feature Selector has been proposed (SU-PSO) towards the selection of optimal feature subset for the improvement of detection performance of classifiers. The feature space ranking is done through the Symmetrical Uncertainty method. Further, memetic operators of PSO include features and remove features are used to choose relevant features and the best of best features are selected using PSO. The proposed feature selector efficiently removes not only irrelevant but also redundant features. Performance metric such as classification accuracy, subset of features selected and running time are used for comparison.

## 1 INTRODUCTION

THE expense of medical and human services is expanding more quickly than the status and the capacity to pay for it. At the same time, more and more data are becoming accumulated due to the availability of computers. Such a huge amount of information cannot be handled by the experts to make diagnosis, prognosis and treatment schedules in short time. Today, the main challenge faced by healthcare industry including hospitals or healthcare centres is delivery of eminence service with affordable cost to the society. Due to misdiagnosis of disease, more people lose their lives. The facility may include correct diagnosis of patients and effective treatments at lower cost. Poor clinical diagnosis might cause deaths. The hospitals must also minimize the cost of medical test by achieving results using the appropriate decision support system and techniques employed. Ultimately, the existing abundant medical data (patient records) raise a query "How to explore the data to get potential idea and make use of this idea for predicting the absence or presence of disease". Patient's healthcare data are stored regularly but are not used to extract useful knowledge. In recent years many research have been conducted to effectively utilize statistical analysis and data mining techniques for diagnosis of disease. Preprocessing is a vital step in the knowledge discovery process. Enhancing the medical database improves the quality of medical diagnosis. The proposed technique focused on pre-processing for data reduction and also for choosing prominent features to improve classification accuracies which help in the prediction of results especially in the field of healthcare.

In data mining, one of the pre-processing steps called Feature selection aims to select relevant set of attributes from the dataset and that would offer better predictive accuracy compared to the model that has been contributed with a complete set of features. The main objective is to decrease the size of the dataset to be prepared by the classifier, enhance the predictive accuracy and furthermore decrease the computational time. Feature selection method involves two general types, the filter and the wrapper approach. Filter method relies upon normal attributes of the preparation information to choose few attributes without including any learning calculation. It surveys the significance of attributes from the information alone, without classifiers and utilizes measures like separation and data consistency. The wrapper method needs one predetermined learning algorithm to assess and figure out which attributes have to be chosen. But, this method has higher computational complexity and requires more time for big dimensional data. But the embedded method incorporates diverse evaluation criteria during various search phases and hence benefits from both the approaches. This approach is

capable of achieving accuracy of a wrapper method at the speed of filter method. The prime goal of this proposed work is to demonstrate that the physician reaches better finding by using the critical features from the complex medicinal dataset. The basic concentration falls on the reduction of dimensionality by selecting meaningful features to improve the predictive accuracy and to assess the proficiency of the SU-PSO Feature Selection strategy. Then the chosen subset is assessed by utilizing four regular classifiers, for example, Naïve Bayes (NB), K-Nearest Neighbor (KNN), Multi Layer Perceptron (MLP) and Support Vector Machine (SVM). The performance metric used are Number of features selected, Accuracy and Running time as execution measurements. The experimental outcomes demonstrate that the proposed SU-PSO based feature selection method accomplishes dimensionality reduction significantly and improves the predictive accuracy with reduced computational time from the datasets acquired from the UCI and Kentridge Bio medical data storage.

## 2    RELATED WORK

MEDICAL data are the real world data and certainly complex and huge in number .Lot of data mining techniques have been proposed in discovering effective techniques for medicinal conclusion. Sellapan et al. (2008) and Asha et al (2010) have structured an Intelligent Heart Disease Prediction model to forecast coronary illness by utilizing three classifiers, such as Neural Networks, Naïve Bayes and Decision Tree. Among these Naive Bayes is evident with great forecast likelihood of 96.6%. Vikas Chaurasi and Saurabh Pal (2013) worked on heart disease process using data mining methods. To predict heart disease, they utilized popular algorithms such as CART, ID3 and decision Table. The performances of the three classifiers were evaluated, and they resulted 83.49%, 72.93% and 82.50% respectively. Liu et.al, (1996) introduced a method based on consistency and it evaluated the significance of features by the strength of consistency in the class values when the training features were proposed onto the subset of attributes. Laetitia et al. (2001) were interested in discovering genetic features and environmental factors that involved in multi- factorial diseases such as obesity and diabetes. Choubey and Sanchita (2016) handled genetic algorithm and multilayer perceptron techniques to diagnose diabetics. This was implemented in two levels. In the first level, feature selection was performed by GA and classification of selected characteristics was implemented using MultiLayer Perceptron Neural Network (MLPNN). This drastically improved the accuracy. Hybrid GA/SVM approach was used by Huerta (2006). Fuzzy logic was used to decrease the size of the initial problem and the redundant genes were eliminated. GA was utilised to extract a subset of better performing genes. Then it was evaluated by SVM. Hsieh et al. (2012) worked on

ensembled machine learning model for diagnosing of breast cancer. Here, feature selection is done using information-gain. The classifiers that were used for developing ensemble classifier were quadratic classifier (QC), neural fuzzy (NF), and the k-nearest neighbor (KNN), The results explained that ensemble framework focused better working performance than single classifier. A hybrid model was developed by Swati et al.(2013) to get better classification performance in the prediction of cardiovascular disease . The algorithms such as Forward Feature Inclusion, Back elimination Feature Selection and Forward Feature Selection were incorporated in the model. By using distance criterion, the features were ranked. Further, the classification of the proposed model was evaluated. These constrain created ranks for all features according to their significant targeting of class identification.Yang and Zhang (2009) have proposed two stages feature selection technique GAEF (Genetic Algorithm with embedded filter). Initially GA is implemented to pre-select features and the filter method is utilised in the next stage for precise sample classification to identify prominent subset of features.

Ding Ding and Peng (2002) have proposed the minimum Redundancy and Maximum Relevance (mRMR) technique which computes Mutual Information (MI). This method computes the correlation of features to find redundant and significant data. Both the MI and the mRMR methods were used to generate feature ranking lists. The aim of this study dealt with the comparison of Artificial Neural Network and Support Vector Machine. The performances of both the methods were compared using BUPA Liver Disorder Dataset. The GA based multilayer perceptron techniques have proposed by Choubey and Sanchita (2016). This method was executed in two levels where, GA and MultiLayer Perceptron Neural Network were used for the diagnosis of diabetics. Olaniyi et al. (2015) introduced a new method for the categorisation of heart disease. Factors for heart diseases and its difficulties with remedies were considered in this work. Hassanien (2004) approached rough set theory to reduce the attribute and provided rule to find breast cancer. Unler et al. (2011) suggested a technique which integrates filter approach with the wrapper methods. The filter method worked based on the mutual information. Moreover in the wrapper model, customized discrete PSO method is used. Initially, they implemented feature selection techniques in medical data to lessen the insignificant attributes. Then, wrapper method was applied to theatrically decrease the cost. A rough set method to generate diagnostic rule was proposed by Tsumoto (2004). This method is based on the hierarchical structure of different medical diagnostics. The generated rules by this method can correctly represent perfect decision process. Gopala Krishna Murthy Nookala et al., (2013) compared 14 different classification algorithms through 3 different types of

cancer data sets. Most of the methods provided enhanced results, when the dimension of the attributes increased. However, accuracy level depended on the kind of the datasets to be used. Finally, they realized that the algorithms do not provide better accuracy level and the user tried to choose the best data set. By using data mining techniques, Cheng-Mei Chen and Chien Yen Hsu (2011) proposed a Survival prediction model for liver cancer. They obtained dataset from the medical data center in North Taiwan during 2004 – 2008. ANN and CART were involved in prediction model. The model was tested with three criteria. And it made a conclusion that ANN model gave more accuracy than the CART model. For high dimensional cancer data a novel feature selection approach has been proposed by Barnali Sahu et al., (2012). Initially, the data is grouped using k–means clustering method and genes are ranked using SNR (signal–to–noise ratio) score. Then, the selected subset of features is optimized using PSO. The proposed method produced excellent accuracy than other methods. MonirulKabi et al., (2011) have presented a new method called as HGAFS (Hybrid Genetic Algorithm for Feature Selection). They employed and embedded a new local search operation to improve the feature selection process. A genetic neuro fuzzy system was proposed by Rawat and Burse (2013). GA used for the selection of features and then it is integrated with Adaptive Neuro-Fuzzy Inference System. The proposed method produced better results.

## 3 PROPOSED METHOD

MULTIPLE Objective Evolutionary Algorithms (MOEAs) are known as search methods from biological world inspired by natural selection which includes survival of the fittest. MOEAs follow the basics of single-objective GAs. GAs initiates with a population of random individuals which are updated during succeeding generations. The crossover and the mutation operators aid in the process to introduce new genetic information in the solutions during each generation. At each successive generation, every individual is evaluated as per a fitness function. Individuals possessing high-fitness values are ranked at the top, where as the individuals with lower fitness values are ranked lower. Such individuals may vanish from the population in consecutive generations. The algorithm proceeds for a pre-determined maximum number of generations which is determined as the stopping criteria. The algorithm can also be terminated when no additional enhancement is observed. The main variation between the single-objective GAs and MOEAs is that, the MOEAs process focuses on the strategies utilized for selection and diversification. Different target developmental calculations obtain their motivation from normal choice and survival of the fittest in the natural world. MOEAs pursue the essentials of single-target PSO. PSO starts with a populace of arbitrary feature (called chromosomes)

that is rectified through progressive generation. The hybrid and change administrators are helpful to introduce new planned structure arrangements at every generation .During each progressive generation, every feature is tested by a wellness work. Feature with high-wellness esteems ranks at the best, where as Feature with low-wellness work esteem is positioned lower and they probably vanish from the populace. The calculation maintains for a pre-decided most extreme number of generation or until the point that no extra enhancement is watched. The Ranking of feature space is executed through SU method.

$$SU\left(X_i, X_J\right) = 2\,\frac{IG(S,X)}{H(X_i) + H(X_j)} \tag{1}$$

where

$$IG(X_i, X_j) = H(X_i) - H(X_i/X_j) \tag{2}$$

$$H(X_i/X_j) = \sum_{xi} P(X_i) \sum_{xi} P(X_i/X_j) log_2 P(X_i/X_j) \tag{3}$$

The heuristic search algorithm proceeds from the empty set of features, and utilises the best-first search along with a halting criterion of 4 consecutive non-improving subsets. The subset which comprises the highest merit is chosen during the search.

- Feature-ranking methods through SU rank from the most relevant features to least relevant. This ranking can be used to discard features.
- The PSO algorithms provide a feature relevance weight to each individual feature through optimization approach. The features Selected by SU can subsequently be optimized through PSO based optimization algorithm.

The two main goals of MOEA are used to guide the search towards identifying the Pareto set. The best-known Pareto set must be a subset of the optimal Pareto set and retains a different set of non-dominated solutions.

- Stage 1: A random populace is instated.
- Stage 2: Objective functions for all targets and constraint are assessed.
- Stage 3: Front ranking of the populace is done dependent on the strength criteria.
- Stage 4: Crowding separation is ascertained (for every ith solution of a specific front, density of solutions in its encompassing is assessed by taking normal separation of two arrangements on its either side along every one of the goal. This normal separation is known as crowding distance).
- Stage 5: Selection is performed utilizing swarmed paired competition determination operator.
- Stage 6: Crossover and transformation administrators are connected to produce a posterity populace.
- Stage 7: Parent and posterity populaces are joined and a non-ruled arranging is finished.
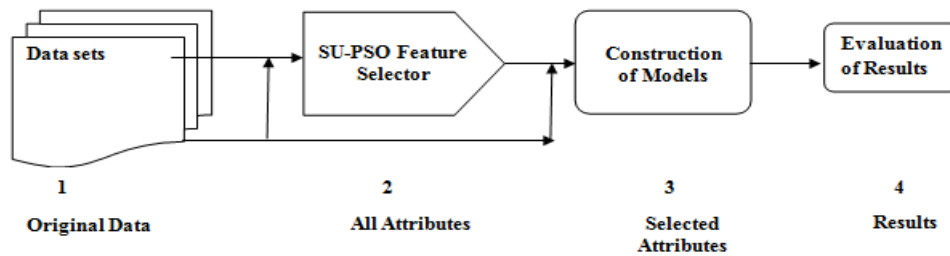
**Figure 1. Proposed SU-PSO Feature Selection Model**

- Stage 8: The best individuals from the joined populace supplant the parent populace.

## 4 RESULTS AND DISCUSSION

THE suggested SU-PSO technique has been implemented on selected medical datasets with the accompanying points of view:
1. Precision of Classifier
2. Number of Feature selected
3. Running time

Name of the Datasets, testing environment, strategy and the vital destinations for the investigations of the assessment of the objectives are portrayed beneath. The proposed approach is executed on various biomedical datasets from the UCI and Kentridge archive. The data consist of little dimensional datasets and medium dimensional datasets and they are shown in Table 1. The speculation capacity of the suggested SU-PSO technique is demonstrated by its process on all these classes of datasets.

Further the proposed calculation creates similarly great characterization precision on all these datasets.

The tests are done on a Windows XP working framework with Intel i5, DDR3, 8GB RAM, and 500GB HDD. The work is actualized in WEKA condition Tool. WEKA has been recognized as a milestone framework in machine learning and information mining and has achieved it's

acknowledgment among the scholarly community and industry. Then, it has been turned into a generally utilized tool for information mining research. Another advantage of its "Open Source" nature is that it offers free access of source code and has empowered to generate and alter the modules for executing the proposed work. The stepwise process is described below. The contribution of the framework is given in the Attribute-Relation File Format (ARFF). Then the calculation is performed and the chosen ideal relevant attributes are fetched as the output. The outcome is done in WEKA by using the name determined in \@relation". The traits determined under \@attribute" and occurrences noted under \@data" are recovered from the ARFF record and added to the table. 10-fold cross validation has been performed for every one of the classifiers. In each run, the dataset has been classified as preparing and testing set, arbitrarily. The outcomes are depicted in Tables 2-4.

The Tables 3 and 4 abbreviate the characterization execution as far as normal exactness and processing time of the predetermined traditional classifiers on the entire data before implementing the SU-PSO feature Selector and the reduced informational collection subsequent to implementing the proposed strategy respectively. On looking at the consequences of traditional classifiers, it is seen that every one of the outcomes is more solid.

**Table 1. Medical datasets used for the experimentation**

| Dataset Name | No. of Instances | No. of Features | No. of Classes |
|---|---|---|---|
| Haberman's Survival | 306 | 4 | 2 |
| Liver Disorder | 345 | 7 | 2 |
| Breast Cancer | 286 | 10 | 2 |
| Cardiac Arrhythmia | 452 | 280 | 16 |

**Table 2.** Features by each feature selection algorithm

| Dataset | All | SU | GA | SU-GA | Proposed SU-PSO |
|---|---|---|---|---|---|
| Haberman's Survival | 4 | 4 | 2 | 2 | 2 |
| Liver Disorder | 7 | 5 | 1 | 1 | 1 |
| Breast Cancer | 10 | 9 | 5 | 5 | 5 |
| Cardiac Arrhythmia | 280 | 103 | 26 | 20 | 14 |

**Table 3a.** Accuracy of SVM and NB classifiers on selected features by each feature selection algorithm
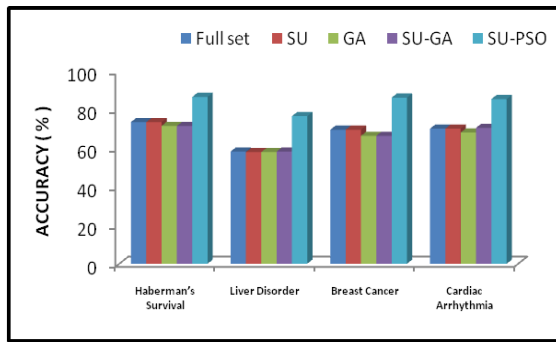
| Dataset | SVM | | | | | NB | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Full set | SU | GA | SU-GA | Proposed SU-PSO | Full set | SU | GA | SU-GA | Proposed SU-PSO |
| Haberman's Survival | 73.52 | 73.52 | 71.56 | 71.56 | **86.63** | 76.14 | 75.16 | 76.12 | 79.86 | **90.65** |
| Liver Disorder | 58.26 | 57.97 | 57.97 | 58.26 | **76.63** | 55.36 | 55.36 | 57.97 | 64.56 | **85.98** |
| Breast Cancer | 69.58 | 69.58 | 66.43 | 66.43 | **86.27** | 71.67 | 71.89 | 72.37 | 80.56 | **85.65** |
| Cardiac Arrhythmia | 70.13 | 70.13 | 68.14 | 70.57 | **85.51** | 62.38 | 66.15 | 74.45 | 78.45 | **85.00** |
| *Average* | 67.87 | 67.80 | 66.03 | 66.71 | **83.76** | 66.39 | 67.14 | 70.23 | 75.86 | **86.82** |

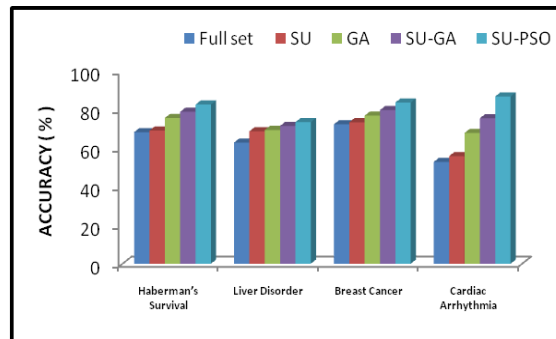**Table 3b.** Accuracy of KNN and MLB classifiers on selected features by each feature selection algorithm

| Dataset | KNN | | | | | MLP | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Full set | SU | GA | SU-GA | Proposed SU-PSO | Full set | SU | GA | SU-GA | Proposed SU-PSO |
| Haberman's Survival | 68.30 | 69.13 | 75.68 | 78.89 | **82.65** | 66.56 | 71.89 | 76.56 | 81.43 | **87.65** |
| Liver Disorder | 62.89 | 68.75 | 69.45 | 71.56 | **73.54** | 57.97 | 74.78 | 79.21 | 85.46 | **88.12** |
| Breast Cancer | 72.37 | 73.48 | 76.89 | 79.78 | **83.73** | 65.23 | 68.45 | 69.56 | 70.12 | **72.76** |
| Cardiac Arrhythmia | 52.87 | 55.89 | 67.89 | 75.56 | **86.86** | 67.25 | 79.42 | 81.56 | 84.15 | **88.51** |
| *Average* | 64.11 | 66.81 | 72.48 | 76.45 | **81.70** | 64.25 | 73.64 | 76.72 | 80.29 | **84.26** |

**Table: 4** Average Accuracy and Running Time of various methods in comparison

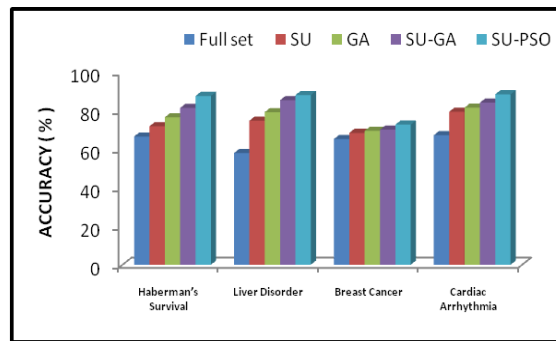| Dataset | Classification Accuracy ( % ) | | | | | Running Time ( Sec ) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Full set | SU | GA | SU-GA | Proposed SU-PSO | Full set | SU | GA | SU-GA | Proposed SU-PSO |
| *Support Vector Machine* | 71.33 | 71.14 | 69.46 | 70.78 | **85.24** | 1.02 | 0.83 | 0.80 | 0.68 | **0.57** |
| *Naïve Bayes* | 70.01 | 70.61 | 73.34 | 78.06 | **87.29** | 0.36 | 0.39 | 0.8 | 0.96 | **0.79** |
| *K-Nearest Neighbour* | 67.41 | 69.56 | 74.29 | 77.67 | **83.26** | 1.27 | 1.06 | 0.90 | 1.18 | **0.92** |
| *Multilayer perceptron* | 71.14 | 78.66 | 81.17 | 83.94 | **87.37** | 2.19 | 2.72 | 2.84 | 4.13 | **3.33** |
| *Average* | 69.97 | 72.49 | 74.57 | 77.61 | **85.79** | **1.21** | **1.25** | **1.34** | **1.74** | **1.40** |

*(a) Support Vector Machine*
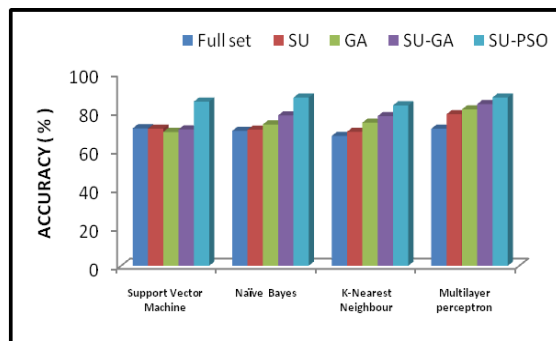


*(b) Naïve Bayes*

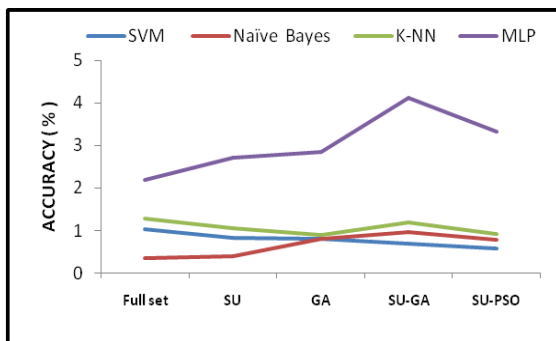

(c) *K-Nearest Neighbour*



(d) *Multilayer perceptron*

**Figure 2.** The Influence of Feature selection on classification accuracy for selected datasets



*(a) Classification Performance*



*(b) Average Running Time of the classifiers*

**Figure 3** Average Performance of the various methods in comparison

The analyzation of the introduced SU-PSO has been performed by well known classifiers by utilizing the determined destinations as the principle measurements. The three principle targets, accuracy of classifier on the selected subset and number of selected features and time taken by the classifier to build the model have been recorded and presented in the Tables 2 to 4 for various dimensional datasets. In Tables 3 and 4, the assessment measures are illustrated both the original data and the features selected by the proposed SU-PSO technique and they are contrasted with different strategies determined by various

authors. From the measures, it is clear that the proposed SU-PSO stands prevalent compared to different techniques.

## 5    CONCLUSION

A novel feature selector has been proposed by incorporating SU and PSO for multi-class grouping by dealing with various dimensional data. The framework is done for making enhancements of the current work with three points of view, for example, decrease in list of original features, increased classification accuracy

and limiting the classifier running time. The proposed SU-PSO which is utilized as the effective feature selector, improves the classifier performance with best precision rate for some high dimensional datasets with least number of features as well as least running time. Likewise, the prevalence of the proposed strategy has been contrasted with three officially existing techniques with the guidance of four traditional classifiers. Additionally, in light of affectability and specificity, it is used in learning calculation and it has been focused on obtaining the main position as well as effectiveness improvement. An examination made on four restorative datasets abridges the qualities of this proposed strategy with different execution measurements namely Classification accuracy, Running Time and Number of features chosen. It is clear from the results that the proposed framework performs well for medical datasets with distinctive number of tests, features and classes, which help in the prediction of diseases especially in the field of healthcare.

## 6    DISCLOUSRE STATEMENT

THE authors report no conflict of interest.

## 7    REFERENCES

Asha Rajkumar and Sophia Reena G, (2010). Diagnosis Of Heart Disease Using Datamining Algorithm, GJCST.

Barnali Sahu and Debahuti Mishra A, (2012). Novel Feature Selection Algorithm using Particle Swarm Optimization for Cancer Microarray Data, International Conference on Modeling Optimization and Computing (ICMOC) 38 pp. 27 - 31.

Blake C.L and Merz C.J. (2008), CI Repository of Machine Learning Databases, http://www.ics.uci.edu/~mlearn /mlrepository.html.

Cheng-Mei Chen, Chien-Yeh Hsu, Cheng-Mei Chen and Chien-YehHsu, (2011). Prediction of survival in Patients with Liver Cancer using Artificial Neural Networks and Classification and Regression Trees, IEEE Seventh International Conference on Natural Computation.

Choubey, and Sanchita D K, (2016). P.: GA_MLP NN: A hybrid intelligent system for diabetes disease diagnosis. Int. J. Intell. Syst. Appl. **8**(1), 49.

Ding, Chris, and Hanchuan Peng, (2002). Minimum Redundancy Feature Selection From Microarray Gene Expression Data." Journal of Bioinformatics and Computational Biology 3.2 (2005): 185-205. 1-6.

Gopala Krishna Murthy Nookala et.al. (2013). Performance Analysis and Evaluation of Different Data Mining Algorithms used for Cancer Classification. International Journal of Advanced Research in Artificial Intelligence, Vol. 2, No.5, pp.49-55.

Hassanian A E, (2004). Rough set approach for attribute reduction and Rule generation: a case of patients with suspected breast cancer, J. Am. Soc. Inform. Sci Technology .55(11) pp.954-962.

Hsieh S L, Hsieh S H, Cheng P H,(2012). Design ensemble machine learning model for breast cancer diagnosis, Journal of Medical Systems, vol. 36, no. 5, pp. 2841–2847.

Huerta E B, (2006). A Hybrid GA / SVM Approach for Gene Selection and Classification of Microarray Data, pp. 34–44.

Kohavi R, John G H, (1997). Wrappers for feature subset selection, Artificial Intelligence.

Laetitia Jourdan, Clarisse Dhaenens & El-Ghazali Talbi A, (2001). Genetic Algorithm for Feature Selection in Data-Mining for Genetics. 4th Metaheuristics International Conference 29-33.

Liu H and Motoda H (1998). Feature Selection for Knowledge Discovery and Data Mining. Boston:Kluwer Academic Publishers, ISBN 0-7923-8198-X.

Liu H, Setiono R A, (1996). Probabilistic approach to feature selection - A filter solution, International Conference on Machine Learning 319–327.

MonirulKabi M D, Shahjahan M D, Murase Kazuyuki, (2011). A new local search based hybrid genetic algorithm for feature selection, Int. J. Neurocomput. 74 (17) 2914–2928.

Olaniyi E. O., Oyedotun O K, and Adnan K, (2015).Heart diseases diagnosis using neural networks arbitration, International Journal of Intelligent Systems and Applications (IJISA), vol. 7,No. 12, p. 75, 2015.

Pablo Bermejo, Luis de la Ossa, Jose A. Gamez, Jose M. Puerta, (2012). Fast wrapper feature subset selection in high-dimensional datasets by means of filter re-ranking, Knowledge- Based Systems.

Rawat K, Burse K, (2013). "A Soft Computing Genetic – Neuro fuzzy Approach for Data Mining and Its Application to Medical Diagnosis," International Journal of Engineering and Advanced Technology (IJEAT), ISSN: 2249 - 8958, no. 1, pp. 409-411, 2013.

Sallehuddin R, Ubaidillah S H and Mustaffa N H, (2014) Classification of Liver Cancer Using Artificial Neural Network and Support Vector Machine", Elsevier Science Proc. Of Int. Conf on Advance in communication Network, and Computing, CNC.

Sellappan Palaniappan and Rafiah Awang, (2008). Intelligent Heart Disease Prediction System Using Data Mining Techniques, IEEE.

Swati S and Ashok G, (2013). Feature selection for medical diagnosis: Evaluation for cardiovascular diseases, Expert systems with applications, vol. 40, pp. 4146–4153.

Tsumoto S, (2004). Mining Diagnostic rules from clinical databases using rough sets and medical diagnostic model, Information Sciences. 162 (2) pp 65-80. 2004

Unler. A, Mura.A and R. B. Chinnam, (2011). Mr 2 PSO: a maximum relevance minimum redundancy approach based on swarm intelligence for support vector machine classification, Inf. Sci., vol. 181, no. 20, pp. 4625–4641.

Vikas Chaurasia, Carib., (2013). Early Prediction of Heart Diseases Using Data Mining Techniques, Caribbean Journal of Science and Technology, SciTech, Vol.1, 208-217, ISSN:0799-3757.

## 8  NOTES ON CONTRIBUTOR

**Ilangovan Sangaiah** is doing his Research in Anna University, India in the area of high dimensional data selection .He is serving in KLN College of Engineering as an Associate Professor of Information Technology. His area of Interest is in Data mining and Wireless Networks. He published papers in several prestigious journals and conferences on this research.

Email:  ilangovans@yahoo.com

**A. Vincent Antony Kumar** is a professor and head of the department of Information Technology, PSNA College of Engineering and Technology, Dindigul, India. He has more than 21 years of teaching experience. His field of Interest includes Cloud computing, sensor network and Data mining. He organized more number of national and International conferences in his credit. He has published papers in several prestigious journals and conferences on this research including in Elsevier and Springer.

Email:  vincypsna@rediffmail.com