# LEAST SQUARES SUPPORT VECTOR MACHINES PARAMETER OPTIMIZATION BASED ON IMPROVED ANT COLONY ALGORITHM FOR HEPATITIS DIAGNOSIS

**Nursuci Putri Husain, Nursanti Novi Arisa, Putri Nur Rahayu, Agus Zainal Arifin, and Darlis Herumurti**

Department of Informatics, Faculty of Information Technology,
Institut Teknologi Sepuluh Nopember (ITS)
Kampus ITS, Surabaya, 60111.

E-mail: nursuci.husain15@mhs.if.its.ac.id, nursanti15@mhs.if.its.ac.id, putri15@mhs.if.its.ac.id

**Abstract**

Many kinds of classification method are able to diagnose a patient who suffered Hepatitis disease. One of classification methods that can be used was Least Squares Support Vector Machines (LSSVM). There are two parameters that very influence to improve the classification accuracy on LSSVM, they are kernel parameter and regularization parameter. Determining the optimal parameters must be considered to obtain a high classification accuracy on LSSVM. This paper proposed an optimization method based on Improved Ant Colony Algorithm (IACA) in determining the optimal parameters of LSSVM for diagnosing Hepatitis disease. IACA create a storage solution to keep the whole route of the ants. The solutions that have been stored were the value of the parameter LSSVM. There are three main stages in this study. Firstly, the dimension of Hepatitis dataset will be reduced by Local Fisher Discriminant Analysis (LFDA). Secondly, search the optimal parameter LSSVM with IACA optimization using the data training, And the last, classify the data testing using optimal parameters of LSSVM. Experimental results have demonstrated that the proposed method produces high accuracy value (93.7%) for the 80-20% training-testing partition.

**Keywords:** *Classification, Least Squares Support Vector Machines, Improved Ant Colony Algorithm, Local Fisher Discriminant Analysis, Hepatitis Disease.*

**Abstrak**

Banyak metode klasifikasi yang mampu mendiagnosa seorang pasien mengidap penyakit Hepatitis, salah satunya adalah menggunakan metode klasifikasi Least Squares Support Vector Machines (LSSVM). Terdapat dua parameter yang sangat berpengaruh pada LSSVM yaitu parameter kernel dan parameter regularisasi. Penentuan parameter optimal tersebut harus diperhatikan untuk mendapatkan akurasi klasifikasi yang tinggi pada LSSVM. Penelitian ini mengusulkan metode optimasi Improved Ant Colony Algorithm (IACA) dalam penentuan parameter optimal LSSVM untuk mendiagnosa penyakit Hepatitis. IACA membuat penyimpanan solusi untuk menjaga rute dari keseluruhan semut. Solusi yang disimpan adalah nilai parameter LSSVM. Ada 3 tahapan utama pada penelitian ini yaitu, dimensi dataset Hepatitis direduksi menggunakan metode Local Fisher Discriminant Analysis (LFDA), kemudian parameter optimal LSSVM dicari dengan metode optimasi IACA menggunakan data training, setelah itu data testing diklasifikasikan menggunakan parameter optimal LSSVM. Hasil uji coba menunjukkan bahwa metode yang diusulkan menghasilkan nilai akurasi yang tinggi (93,7%) pada partisi 80-20% training dan testing.

**Kata Kunci:** *Klasifikasi, Least Squares Support Vector Machines, Improved Ant Colony Algorithm, Local Fisher Discriminant Analysis, Hepatitis.*

## 1. Introduction

LSSVM classification method is proposed by Suykens, et al. [1], and LSSVM is a development method of the SVM [2]. In the SVM, the optimal hyperplane is obtained by solving quadratic programming problem by minimizing a function with an inequality condition. Different with SVM, LSSVM gives solutions with linear equations, not with the quadratic programming problems [3]. There are two parameters that very influence to improve the classification accuracy on LSSVM, they are kernel parameter ($\sigma^2$) and regularization parameter ($\gamma$). Determining the optimal parameters must be considered to obtain a high classification accuracy on LSSVM. The parameters can be searched by trial and error, but trial and error is very not efficient and not effective.
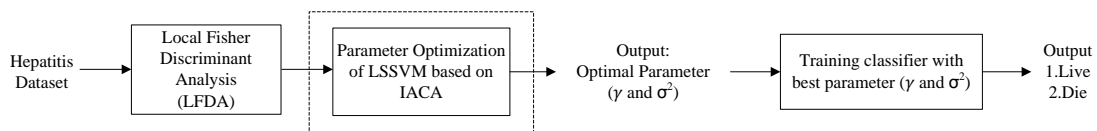
Figure 1. Hepatitis classification system

Therefore, many studies used the Cross Validation (CV) to optimize the parameters LSSVM. However, CV has an disadvantage in computing speed and the accuracy of classification still in average. So that, many researchers proposed a method for optimizing the parameters of the LSSVM.

Zhigang and Chengling [4] proposed a method for predicting the damage depth of coal seam floor using LSSVM. The optimal parameters of LSSVM optimized by Particle Swarm Optimization (PSO). PSO is one of the optimization methods inspired by the behavior of a group of animal movements such as the movement of a group of birds (flock). Each object of animals becomes a particle. A particle in search space has a position that encoded as vector coordinates. This position vector is considered as a state of being occupied by a particle in the search space. Each position in the search space is an alternative solution that can be evaluated using the objective function. Based on experimental results for their testing and training data, PSO-LSSVM can measure the depth of the coal seam floor damage.

Gupta et al [5], proposed a hybrid method of Genetic Algorithm (GA) and LSSVM to increase fault classification of the power transformer. GA generate the initial population randomly, expanding the search space, improve the speed of convergence and search the optimal parameters. The GA-SVM method successfully improve the accuracy of classification errors on the power transformer using DGA dataset (Dissolved gas analysis).

Then, Hegazy et al. [6], proposed a hybrid method Artificial Bee Colony (ABC) and LSSVM on stock price prediction. ABC algorithm is an algorithm inspired by the habits of bees exploration (foraging) to find the optimal solution. ABC chose the best parameters for the LSSVM and avoid the over-fitting problem. That study compared the proposed method with PSO optimization method, where the method of ABC-LSSVM has high convergence speed than PSO-LSSVM. However, the ABC method produces a local minimum parameter [7].

This study proposed an optimization method based on Improved Ant Colony Algorithm (IACA) in determining the optimal parameters in LSSVM to diagnose Hepatitis disease. This algorithm aims to find the optimal path. The search path is based on the behavior of ant colonies in finding the path to a food source [8]. This basic idea then used to solve the problems which illustrated by the behavior of ants. IACA is an optimization method that makes the storage solution to keep the whole route of the ants. The solutions that have been stored were the value of the parameter LSSVM. IACA produces a global minimum parameter at the end of the iteration.

The remainder of this paper is organized as follows. Section 2 describes the methods. Section 3 explained the experimental results. And finally, conclusions and recommendations for future work are summarized in Section 4.

## 2. Methods

### Local Fisher Discriminant Analysis (LFDA)

The dataset that have large feature dimension can be affected to the classification process. The feature dimension can be reduced with dimensionality reduction method. According to [9], dimensionality reduction method is divided into two, they are feature extraction and feature selection. Feature extraction is one of dimensionality reduction method to looking for features that have most relevant information to the original data by transforming the input data into a set of data with the feature that have been reduced [10]. There are several stages in this study (Figure 1).

This study used feature extraction method called Local Fisher Discriminant Analysis (LFDA) to reducing the feature dimension of the dataset. LFDA proposed by Sugiyama [11], LFDA maximize between class separation and defending within class local structure [10]. $S^{(bc)}$ and $S^{(wc)}$ are scattered matrix of between class and within class, both of them calculated by equation (1) and (2) [10].

$$S^{(bc)} = \frac{1}{2} \sum_{ij=1}^{n'} W_{i,j}^{(bc)} (x_i - x_j)(x_i - x_j)^T, \tag{1}$$

$$S^{(wc)} = \frac{1}{2} \sum_{ij=1}^{n'} W_{i,j}^{(wc)} (x_i - x_j)(x_i - x_j)^T, \tag{2}$$

**Algorithm 1:** Parameter Optimization of LSSVM based on IACA

Input:  Number of solutions (*N*), number of ants (*m*), range parameter value ($\gamma$, $\sigma^2$), size of solutions storage (*k*), termination criterion

Output: Optimal parameter values ($\gamma$, $\sigma^2$) for LSSVM and classification accuracy

Begin

    Initialize *N* solutions

    Call LSSVM to evaluate *N* solutions

    //Sort solutions and save them in solutions storage

    *A* = Sort(*S0,S1,.....Sn*)

    While termination criterion is not do

    //Generate m solutions

    For *i* = 1 to *m* do

        //Build solution

        Choose *S* according to its weight vector

        Save new solution

        Call LSSVM to evaluate new solutions

    End

    // Sort solutions and choose the best *N*

    *A* = best(Sorting *S0, S1, ..Sn + m*), *N*)

End

$n$ is the number of sample in the dataset while $(x_i - x_j)$ is the value based on the local scaling approach [11]. Then the transformation matrix of LFDA $T^{(M)}$ defined in equation (3) [10].

$$T^{(M)} = \frac{\arg max}{T \epsilon \mathbb{R}^{dxr}} \lfloor tr(T^T S^{(bc)} T (T^T S^{(wc)} T)^{-1} \rfloor, \qquad (3)$$

$d$ is the number of dataset dimension, and $r$ is the feature dimension that have been reduced. LFDA search the transformation matrix $T$ of the scatter space between class $T^T S^{(bc)} T$ the distribution are maximized and the scatter space within class $T^T S^{(wc)} T$ the distribution are minimized. Dataset dimension that have been reduced is divided into two training – testing partitions, namely 70-30% and 80-20%.

**Parameter Optimization of Least Squares Support Vector Machines based on Improved Ant Colony Algorithm**

*Least Squares Support Vector Machines*
Least Squares - Support Vectors Machines (LSSVM) is one of modification of SVM [2] that have been proposed by Suykens and Vandewalle

[1]. Besides, the complexity of the calculation is lower, training process LSSVM in large scale is also faster and computation resource is lower than SVM. The same as SVM, LSSVM can be used to classification problem and regression both in the linear case or nonlinear. In the nonlinear case, kernel technique can be applied in the LSSVM. Kernel option that can be used the same as SVM, they are linear, polynomial, RBF, and MLP [1].

$$Q(w, b, \alpha, \xi) = \frac{1}{2} \|w\|^2 + \frac{C}{2} \sum_{i=1}^{n} \xi_i^2$$
$$- \sum_{i=1}^{n} \alpha_i \{y_i[(w.x_i) + w_0] - 1 + \xi_i\} \qquad (4)$$

Every training set is expressed by $(x_i. y_i)$ with $i = 1, 2, \ldots N$, then $x_i = (x_{i1}, x_{i2}, \ldots, x_{iq})$ are attribute or feature for the training set $i$, $y_i = (-1, +1)$ the class label. Then $\|w\|$ is the weight vector of $w$, $C$ is a parameter that used to control the trade-off between margin and classification error. Then, $n$ is the number of data, and $\xi$ is the slack variable. LSSVM trained by minimizing two function equation using Lagrange Multiplier and become equation (4) [1].

The next difference between SVM and LSSVM are $\alpha_i$ (Lagrange multipliers), the value $\alpha_i$ in LSSVM is positive or negative, whereas in SVM the value must be positive. Besides that, if using the RBF kernel, the number of parameters that must be optimized on the LSSVM lower than SVM, in LSSVM only two parameters that need to determined, they are kernel parameter ($\sigma^2$) and regularization parameter ($\gamma$). These parameters can be optimized to make the LSSVM hyperplane separating the classes optimally, even in the high dimension space.

*Improved Ant Colony Algorithm Optimization*
Ant colony algorithm proposed by Marco Dorigo [8] is a probabilistic technique used to solve optimization problems of computing with finding the optimal value to a parameter. The optimal value is the value that obtained through a process and considered to be the best solution of all existing solutions. Ant Colony Algorithm deal with discrete and continuous functions. However, Ant Colony Algorithm that deals with continous functions is considered as a research field [8].

Improved Ant Colony Algorithm (IACA) is modification algorithms based on ant colony algorithm. This algorithm is applied to seen clearly the optimization of continuous functions with increasing several algorithms, such as the objective function below:

$$\min f(x_1, \ldots x_n), x_i \in [a_1, b_1], i = 1, 2, \ldots, n \tag{5}$$

Firstly, initializing the number of solution N, then define the range of parameters ($\gamma$ and $\sigma^2$), m shows the population of ants, $x_i^{(0)} = (x_1, x_2, \ldots x_n)$ shows the initial position toward the destination position and $x_i^{(0)}$ is a random point in range variable *i*.

IACA create a storage solution to keep the route of the overall ants. The solution that have been saved is the parameter value ($\gamma$ and $\sigma^2$) LSSVM. The storage of solution need transition probability equation called the weight vector (*w*), *w* will calculate the solution that have been saved in storage solution with the following equation:

$$w_t = \frac{1}{Qk\sqrt{2\pi}} e^{\frac{(t-1)^2}{2Q^2k^2}} \tag{6}$$

Q is the parameter that controls the process of finding the solution, and k is the size of the storage solution. Our proposed algorithm can be seen in Algorithm 1.

Classification accuracy in this algorithm is used to update the storage solution. Then, the transition probability equation is used to choose the solution route of an ant. The solution of the ant will be used as parameters ($\gamma$ and $\sigma^2$) in the kernel RBF from LSSVM classification.

**Training LSSVM using the optimal parameter**

After obtained the classification model, the next step conducts the prediction process on testing data. In this study, kernel RBF is used to LSSVM classification because its ability to handled the high

dimension data [2] and produce a good performance [12].

3. **Results and Analysis**

In order to evaluate the effectiveness of the proposed method, this study conducts experiments on the Hepatitis dataset. Hepatitis dataset that have been used in this study is from the KEEL Repository [13]. The aim of this dataset is to predict whether a patient Hepatitis disease will die or still live. Hepatitis dataset consists of 19 attributes, 80 instances, and two class labels, they are "die" or "live". There are 13 class instances labeled "die" and 67 class labeled "live". Hepatitis dataset of KEEL Repository did not contain the missing value. Table 1 is a Table 1nformation 19 attributes of dataset Hepatitis.

Feature extraction method that has been used is LFDA, LFDA algorithm was implemented in Matlab [11]. The number of features extracted using LFDA that we get in the experiment was 5 features. The scatter Plot of class after the dataset dimension reduced can be seen in Figure 2. The sign x indicates the class "live" and o shows the class "die". After getting the reduced dataset dimension, the dataset will be used as input into the next process. The next process is training LSSVM using the Improved ant colony optimization as the system diagram shown in figure 1.

Dataset dimension that have been reduced are divided by two training - testing partitions namely 70-30% and 80-20% as shown in Table 1I. After that, this study makes a classification model and search optimal parameters of LSSVM based on IACA optimization using the training data.

The Input parameters that have been used in IACA process in search the optimal parameters of LSSVM as shown in Algorithm 1 are the number

TABEL 1
HEPATITIS DATASET

| No | Attribute | Value |
|---|---|---|
| 1 | Age | 10, 20, 30, 40, 50, 60, 70, 80 |
| 2 | Sex | Male, Female |
| 3 | Steroid | No, Yes |
| 4 | Antivirals | No, Yes |
| 5 | Fatigue | No, Yes |
| 6 | Malaise | No, Yes |
| 7 | Anorexia | No, Yes |
| 8 | Liver Big | No, Yes |
| 9 | Liver Firm | No, Yes |
| 10 | Spleen Palpable | No, Yes |
| 11 | Spiders | No, Yes |
| 12 | Ascites | No, Yes |
| 13 | Varices | No, Yes |
| 14 | Bilirubin | 0.39, 0.80, 1.20, 2.00, 3.00, 4.00 |
| 15 | Alk Phosphate | 33, 80, 120, 160, 200, 250 |
| 16 | Sgot | 13, 100, 200, 300, 400, 500 |
| 17 | Albumin | 2.1, 3.0, 3.8, 4.5, 5.0, 6.0 |
| 18 | Protime | 10, 20, 30, 40, 50, 60, 70, 80, 90 |
| 19 | Histology | No, Yes |

TABLE 2
TRAINING SET AND TESTING SET

| Training – testing partition (%) | Number of instances | |
|---|---|---|
| | Training set | Testing set |
| 70-30 | 56 | 24 |
| 80-20 | 64 | 16 |

TABLE 3
OPTIMAL PARAMETER FOR EACH PARTITION FOUND BY IACA

| Partition ( %) | $\sigma^2$ | $\gamma$ |
|---|---|---|
| 70-30 | 17.3 | 99.8 |
| 80-20 | 18.1 | 100.4 |

TABLE 4
CONFUSION MATRIX FOR EACH PARTITION

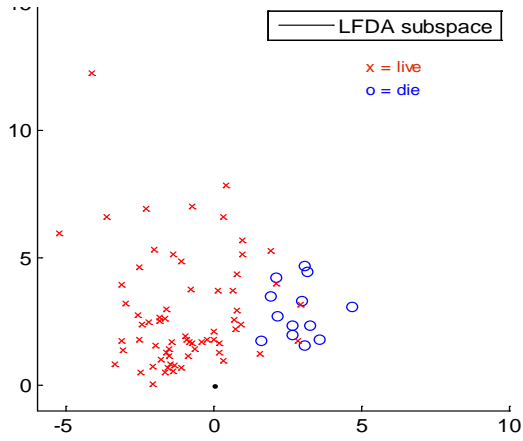| Actual | Predicted | | Partition (%) |
|---|---|---|---|
| | Die | Live | |
| Die | 5 | 0 | 70-30 training-testing |
| Live | 2 | 17 | |
| Die | 3 | 0 | 80-20 training-testing |
| Live | 1 | 12 | |

Figure 2. scatter plot class of the reduced dataset found by LFDA

of solution ($N$) = 50, the number of ants ($m$) = 10, range $\gamma \in$ (0, 150) and range $\sigma^2 \in$ (0, 20).

The optimal parameters obtained by IACA optimization shows in Table 1II. Each partition generate different value of $\gamma$ and $\sigma^2$, for 70-30% training-testing partition gain $\gamma$ = 17.3 and $\sigma^2$ = 99.8. And 80-30% training-testing partition gain $\gamma$ = 18.1 and $\sigma^2$ = 100.4. After classification model obtained, this study conduct prediction process on the testing data.

In order to evaluate the prediction performance of our proposed method, this study computes classification accuracy, sensitivity, and specificity, as shown in confusion matrix for each partition. A classification system is expected to classify all data sets correctly. Generally, the way to measure the performance of a classification using confusion matrix. The confusion matrix is a table that records the result of classification. Based on confusion matrix, it can be seen the amount of data of each class that predicted correctly. By knowing the amount of data that classified correctly, so that it is easy to know the accuracy of the prediction. Another quantity that can be used as a performance classification metric is the sensitivity and specificity. Both of these quantities provide a more relevant performance value. Sensitivity or true positive rate is used to measure the proportions of the original positives correctly predicted as positive. While the specificity used to measure the proportions of the original negatives correctly predicted as negative.

Formula accuracy, sensitivity, and specificity can be seen in equation (7) - (9). In that equation, *TP* (True Positive) is the number of data that is identified properly, *TN* (True Negative) is the number of data that is rejected correctly, *FP* (False Positive) is the number of data that is identified wrongly, and the *FN* (False Negative) is the number of data wrongly rejected.

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN} \; x \; 100\% \qquad (7)$$

$$Sensitivity = \frac{TP}{TP + FN} \; x \; 100\% \qquad (8)$$

$$Specificity = \frac{TN}{FP + TN} \; x \; 100\% \qquad (9)$$

Classification results can be seen using the confusion matrix in Table 1V. We can see in that table the number of false positives more decrease if the size of the training data improved. Then, the classification accuracy using the proposed method can be seen in table V, highest accuracy at 80-20% training - testing partition is 93.7%. The sensitivity that had been achieved for both partitions is 100% while the specificity values for 70-30% training-testing partition was 89% and specificity for 80-20% training - testing partition is 92%.

This study compared the proposed method to classify Hepatitis dataset using LFDA method as a dimensionality reduction method and LSSVM as a classification method without using IACA optimization. Classification accuracy that had been obtained as shown in table VI is 83.3% for the 70-30% training-testing partition and 87.5% for the 80-20% training - testing partition.

This study also classified Hepatitis dataset without using the method of dimension reduction and optimization in determining the optimal parameter and only used the LSSVM classification method. We can see in table VII the classification accuracy that obtained was 79.1% for the 70-30% training-testing partition and 81.25% for the 80-20% training-testing partition.

We also compared the proposed method with previous studies that proposed a classification method for diagnosed Hepatitis disease. Can be seen in table VIII, Genetic Algorithm (GA) and SVM or GA_SVM method [14] proposed by Tan et al., achieve accuracy value as much as 90%. GA is a heuristic search algorithm based on the biological evolution mechanisms. In that study, the GA method used to select the best attributes of

TABLE 5
ACCURACY, SENSITIVITY, SPECIFICITY USING IACA-LSSVM

| Metrics | 70-30% training – testing | 80-20% training - testing |
|---|---|---|
| Accuracy | 91.6 | 93.7 |
| Sensitivity | 100 | 100 |
| Specificity | 89 | 92 |

TABLE 6
CLASSIFICATION ACCURACY USING DIMENSION REDUCTION (LFDA) AND LSSVM WITHOUT IACA OPTIMIZATION

| Partition (%) | Accuracy (%) |
|---|---|
| 70-30 | 83.3 |
| 80-20 | 87.5 |

TABLE 7
CLASSIFICATION ACCURACY ONLY USING LSSVM

| Partition (%) | Accuracy (%) |
|---|---|
| 70-30 | 79.1 |
| 80-20 | 81.2 |

TABEL 8
CLASSIFICATION ACCURACIES OBTAINED WITH OUR METHOD AND OTHER METHODS

| Method | Accuracy (%) |
|---|---|
| CSFNN | 90 |
| LDA | 86,4 |
| GA-SVM | 89,6 |
| Our Method | 93,7 |

Hepatitis dataset, then the dataset that have been selected were classified using SVM with 20 fold cross validation.

While the Local Discriminant Analysis (LDA) method [15] proposed by Stern and got 86,4% of accuracy. In that study, Linear Discriminant Analysis (LDA) modified by Bayes algorithm function called maximum likelihood. LDA is a classification method that tries to find a linear subspace and maximize the separation of two classes based on Fisher Criterion.

And the method proposed by Ozyilmaz and Yildirim namely Conic Section Function Neural Network (CSFNN) [16] achieve an accuracy values as much as 77,4%. The CSFNN is NN algorithm that combined Multilayer Perceptron (MLP) and Radial Basis Function (RBF) to improved the Back Propagation performance.

## 4. Conclusion

This study proposed an optimization method based on Improved Ant Colony Algorithm (IACA) for LSSVM in determining the optimal parameters for diagnosing Hepatitis disease. IACA Algorithm gives optimal parameter LSSVM in each iteration. This study has three main steps: 1) the dimension of Hepatitis dataset reduced by LFDA, 2) search the optimal parameter LSSVM with IACA optimization using the data training, and 3) classify the data testing using optimal parameters of LSSVM. The experimental results show that the proposed method is able to improve the accuracy classification of Hepatitis disease.

This study compared the performance of our method with three other methods, they are LDA, CSFNN, and GA-SVM. Our proposed method achieved high accuracy for the 80-20% training-testing partition (93.7%).

Future investigation will pay attention about the influence of range value $\gamma$ and $\sigma^2$ that we used in search the optimal parameter of LSSVM. Then, Analyzing the input parameter of IACA should be our future work.

## References

[1] Suykens, J. A. K., & Vandewalle, J. "Least squares support vector machine classifiers", *Neural Processing Letters*, vol. 9 (3), pp. 293–300, 1999.

[2] Vapnik, V. "The nature of statistical learning theory". New York: Springer, 1995.

[3] Tsujinishi, D., & Abe, S. "Fuzzy least squares support vector machines for multi-class problems", *Neural Networks Field*, vol. 16, pp. 785–792, 2003.

[4] Zhigang, Yan., Chengling, Cui., "An intelligent model for predicting the damage depth of coal seam floor based on LSSVM optimized by PSO", *Jurnal of applied sciences* 13 (11), pp. 1954-1959, 2013.

[5] Gupta, Aparna R. Et al., "LSSVM Parameter Optimization Using Genetic Algorithm To Improve Fault Classification Of Power Transformer, Engineering Research and Applications", IJERA, Vol.2, Issue 4, pp.1806-1809, July-August 2012.

[6] Hegazy, Osman., Omar S. Soliman, and Mustafa Abdul Salam, "LSSVM-ABC Algorithm for Stock Price prediction", *International Journal of Computer Trends and Technology* (IJCTT) – vol: 7 number 2, Jan 2014.

[7] Gao, W., & Liu, S. (2012). "A modified artificial bee colony", *Computers & Operations Research*, vol: 39, pp. 687-697, 2012.

[8] Dorigo, M. and Stutzle, T., "Ant Colony Optimization", *The Massachusets Institut of Technology Press*, Cambridge, 2004.

[9] Pudil, P.; Novovicová, J. "Novel Methods for Feature Subset Selection with Respect to Problem Knowledge". In Liu, Huan; Motoda, Hiroshi. *Feature Extraction, Construction, and Selection*. p. 101, 1998.

[10] Chen, Hui-Ling. "A new hybrid method based on local fisher discriminant analysis and support vector machine for Hepatitis disease diagnosis", *Internasional Journal of Engineering and science*, vol: 38, pp. 11796-11803, 2011.

[11] Sugiyama, M. "Dimensionality Reduction of Multimodal Labeled Data by Local Fisher Discriminant Analysis", *Journal of Machine Learning Research*, vol. 8, pp.1027-1061, 2007.

[12] Zhang, H. et al., "Three-Class Classification Models of LogS and LogP Derived by Using GA – CG – SVM Approach", *Molecular Diversity*, Springer, vol. 13, no. 2, pp. 261-268, 2009.

[13] J. Alcalá-Fdez, A. Fernandez, J. Luengo, J. Derrac, S. García, L. Sánchez, F. Herrera. KEEL Data-Mining Software Tool: "Data Set Repository, Integration of Algorithms and Experimental Analysis Framework". *Journal of Multiple-Valued Logic and Soft Computing* 17: 2-3, pp. 255-287, 2011.

[14] K.C. Tan, E.J. Teoh, Q. Yu, K.C. Goh, "A hybrid evolutionary algorithm for attribute selection in data mining", *Expert Systems with Applications*, vol: 36 (4), pp. 8616–8630, 2009.

[15] B. Ster, A. Dobnikar, "Neural Networks in Medical Diagnosis: Comparison with Other Methods", 1996.

[16] L. Ozyilmaz, T. Yildirim, "Artificial neural networks for diagnosis of Hepatitis disease", in: *Proceedings of the International Joint Conference on Neural Networks*, 2003, vol. 1, pp. 586–589, 2003.