# Particle Swarm Optimization Feature Selection for Classification of Survival Analysis in Lymphoma Cancer

Norshafarina binti Omar

Faculty of Computing
Universiti Teknologi Malaysia
Skudai, Johor, Malaysia
norshafarina.omar@gmail.com

Mohd Shahizan bin Othman

Faculty of Computing
Universiti Teknologi Malaysia
Skudai, Johor, Malaysia
shahizan@fsksm.utm.my

Roliana binti Ibrahim

Faculty of Computing
Universiti Teknologi Malaysia
Skudai, Johor, Malaysia
roliana@utm.my

Fatimatufaridah binti Jusoh

Faculty of Computing
Universiti Teknologi Malaysia
Skudai, Johor, Malaysia
efaridah88@gmail.com

*Abstract*— **It could be seen that almost survival analysis in biomedical and healthcare are focusing on survival time; related to how long individuals with disease will survive or dying. While most existing survival analysis aims at improving the survival rate by extracting the useful knowledge from patient data, either through existing techniques or through development of new techniques, this paper focused on selecting only the significant and relevant patient data with minimal information loss to classify the patient survival. This paper highlights and discusses the concept and limitation of classification of survival analysis in lymphoma cancer and the abilities of feature selection to solve the classification problems. Therefore, the aim of this paper is to propose particle swarm optimization (PSO) feature selection for the classification of survival analysis in lymphoma cancer. Experiment result from the proposed approach is then compared with the classification without feature selection. From the comparison results, classification of survival analysis with PSO feature selection outperformed classification of survival analysis without feature selection; with 85.45% compared to 77.77% each.**

Keywords — **PSO, feature selection, classification, survival analysis, DLBCL cancer, SVM**

## I. INTRODUCTION

Survival analysis has not been widely adopted and fully explored particularly in lymphoma cases. In recent online search of the IEEE digital library, back from 2004 until 2012, we obtained about 1187 entries with the keyword *survival analysis*. However, with the keyword *cancer survival analysis*, we found about 156 entries and obtained 8 entries with the keywords *lymphoma survival analysis* (only makes up 0.7% of total survival analysis entries). Each researcher has a different definition and standard for survival analysis. Survival analysis brings several meaning which are (1) collection of statistical procedures which accommodate time to event censored (incomplete) data so that reliable and accurate information can be obtained (Liu, 2010; Zhongxin, 2011), (2) failure time analysis or time-to-event analysis which are not necessarily associated with failure or survival at all (Zhanshan and Survival, 2008), (3) analysis of data that

1

corresponds to the survival time or failure time (Hamdan and Garibaldi, 2012). (Zhanshan and Survival, 2008) listed lifetimes of organisms, durations of economic recessions, failure time of machine components, and survival times of cancer patients as examples of survival analysis cases. Despite the prototypical event related to medical cases, (Fox, 2006) listed different events for representing survival analysis. Divorce, child-bearing, unemployment, criminal recidivism and graduation from school can be considered as some events that can be related to survival analysis. In this paper, survival is considered as any incidence of lymphoma where the person is alive or dead from the time of diagnosis depending on the individuals' condition. Survival data for survival analysis usually involves several features. These features are playing significant role in order to generate the survival analysis result. Normally, the selections of parameters are depending on the types of cancer and also the availability of patient's information. For example, some registry does not have full dates (day, month, and year) of diagnosis or does not provide the histological information. However, some survival analysis studies will use the same parameters (age, stage of diagnosis and follow-up status). It is not a new phenomenon in health care industry that we often flooded by cancer data but still lack of useful information since it is obviously data could not tell us anything without processing (Hasan and Tahir, 2010). This motivates the need for sufficient method that capable to extract the useful knowledge.

This paper is organized as follows: In Section I, we present the introduction of this paper. In Section II, we present a study on classification including the abilities and limitations of classification in survival analysis. Then Section III, we discuss the abilities of feature selection in classification problem. Section IV provided the proposed approach framework. Experimental results are presented to prove the proposed approach in Section V. Lastly, Section VI discusses the idea of feature selection for solving classification problem in cancer survival analysis.

## II. CLASSIFICATION

In order to find the underlying patterns and knowledge of the Diffuse Large B-Cell Lymphoma (DLBCL) dataset and to make use of the found patterns and knowledge, we carried out some experiment on DLBCL dataset. With the focus on classifying the patient survival status, SVM classifier is powerful tool for supervised learning and widely used in classification problems (Kumar and Gopal, 2010).

Support Vector Machine classifier is a machine learning technique, who originally introduced by (Corter and Vapnik, 1995). Figure 1 shows the overview of the classification process (Othman, 2008). In summary, SVM works by dividing a dataset into training and testing dataset. Both training and testing dataset are scaling so that training and testing will be faster. The advantages of scaling in SVM are to avoid attributes in greater numeric ranges dominating those in smaller numeric ranges and to avoid numerical

difficulties during calculation. The training dataset is run along with selected parameter setting to build a train model. Then the testing dataset is run by learning from the train model to produce an output.
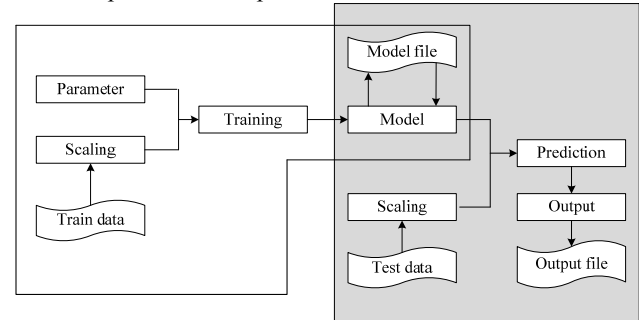


Fig.1. Overview of overall classification process using SVM

However, an issue that could be highlighted from the study of SVM classification is the limitations of classification itself in survival analysis. Some limitations of classification in survival analysis could be stated as follow:

i) In some situations, the proposed classifier is not good enough and do not work well for data which have many features.

ii) Too many features that the classifier is dealing affect the effectiveness of the classifier wherefore many of which will be redundant for the task of classification.

iii) Also too many features that goes through the classifier caused the classifier to work more (workload happen) as well caused the decreasing accuracy of the classification.

## III. FEATURE SELECTION

Based on the limitations of classification as stated earlier in the previous section, feature selection technique is proposed to overcome the drawback of the classifier. Classification requires careful consideration when it comes to dataset before giving the data to classifier. It is better to consider only necessary features rather than adding many irrelevant features since it will makes classification process much harder. So, it is very helpful to have sufficient techniques that capable of selecting the relevant and significant features. Moreover, if feature selection is adopted in classification, it helps in finding the significant feature as well as reduced the workload of the classifier which also improved the classification accuracy. Based on the review of the existing literature (Omar et al, 2013; Jensen, 2005; Rahman et al, 2009; Tu et al, 2007; Sharkawy et al, 2011; Geetha et al, 2008; Liu et al, 2011; Mishra and Sahu, 2011; Ahmed, 2005; Wang et al, 2007), it could be seen that particle swarm optimization enjoy better selection in term of classification accuracy compared to other existing feature selection techniques. From the review (Wang et al, 2007; Blackwell, 2005; Shahamatnia and Ebadzadeh, 2011; Wei et al, 2008; Jamian et al, 2012; Liu et al, 2011; Fan and Wan,

2008; Kennedy and Spears, 1998), some abilities of PSO for features selection could be stated as follow:

i) PSO has powerful exploration ability until the optimal solution is found due to the fact that different particles have the possibility to explore different parts of solution space.
ii) PSO is particularly attractive for feature selection since the particle swarm has memory and knowledge of the solution is retained by all particles as they fly within the problem space.
iii) The other attractiveness of PSO is due to its computationally inexpensive implementation and yet gives decent performance.
iv) The fact is, PSO works with the population of potential solution rather that with single solution.
v) PSO is able to deal with binary and discrete data.
vi) PSO has better performance compared to other feature selection techniques in terms of memory, run time and do not need complex mathematical operators.
vii) PSO is easy to implement with few parameters and is easier to realize and also gives promising results.
viii) The performance of PSO nearly is not affected by the          dimension of the problem.

Particle swarm optimization is a computation technique inspired by simulation of social behavior and proposed by (Kennedy and Eberhart, 2001). The original concept of PSO is to simulate the behavior of flying birds and their means of information exchange to solve problem. The overview to gain insight how PSO works in searching for the significant features is defined in Figure 2 (Wang et al, 2007).
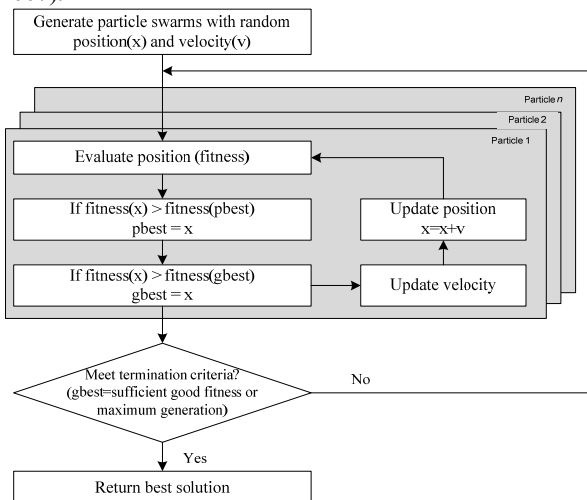


Fig.2. PSO Process

Each particle represent as binary bit of length N, where N is the total number of features. Every bit represents a feature, with the value '1' means the corresponding attribute is selected and '0' not selected. The velocity of each particle is implying how many of the particle's features (bits) should be changed, at a particular moment in time, to be the same as that of the global best position. Here, the velocity of the particle flying toward the best positions. The difference between current particle and best position is indicated by the number of different features (bits) between two particles. Let say, compared with the best position, the current features should be selected but is not, which will lead to a lower fitness values and decreased the quality. Contrary, compared with the best position, the current feature should not be selected, but is selected, and this will lead to redundant features and lead to a lower fitness values too. After particles update their velocity, their position will be updated too. Only by then the particle moves toward the global best. After the particles reaches the global best position, it still exploring the search space for further search. The particles will stop after it found the optimal solution with highest fitness value.

## IV. THE PROPOSED FEATURE SELECTION FOR CLASSIFICATION OF SURVIVAL ANALYSIS APPROACH

Based on the observations made in previous section, we present the idea of PSO feature selection for the classification of survival analysis in this section. The central idea of this approach is to find the optimal number of features which also reduced the workload of classifier as well as improved the classification accuracy of survival analysis. To make it clear, the proposed approach framework is designed in a simple figure. Figure 3 showed the structure of proposed approach, which consists of a data preparation, feature selection, and classification of survival analysis. Each phases contributed an output as each output from previous phase will then lead to the next phases.
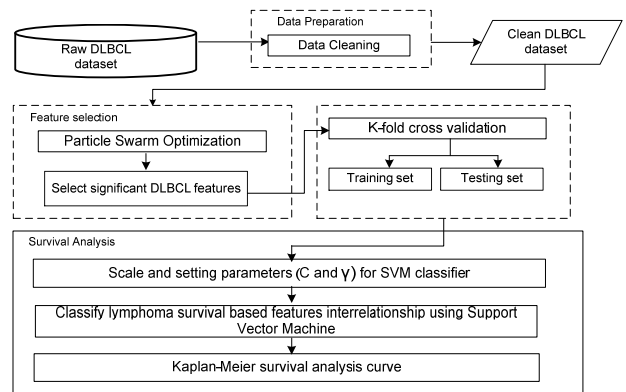


Fig.3. Proposed Approach

The dataset used to find the best significant features is Diffuse Large B-Cell Lymphoma (DLBCL) dataset which is retrieved from The Lymphoma/Leukemia Molecular

Profiling Project. The DLBCL dataset was actually contained of 9 features including the class feature (patient's follow up years, patient's status at follow-up, subgroup, international prognostic index (IPI) group, germinal center B cell signature, lymph node signature, proliferation signature, bone morphogenetic protein 6 (BMP6) and major histocompability complex (MHC class II)) with the total amount of sample is 240. The class feature is to classify each sample of data whether the patient with DLBCL cancer is survived or died.

As shown in Figure 3, the proposed PSO-SVM survival classification is started with data preparation. Data preparation is needed to prepare data for analysis. The process of data preparation would affect the quality and accuracy of the survival analysis. After done with the data preparation, PSO feature selection was run to identify and select the most significant features in DLBCL dataset. Different parameters setting were set up before the PSO is applied. The selected features were chosen based on best features found with highest fitness value.

Next, using the selected features brought from the previous process, the DLBCL dataset with selected data features need to split into training and testing set using *k*-fold cross validation. After the data were split according to *k*-fold cross validation, experiments were conducted with respect to different SVM parameter values of cost (C) and gamma ($\gamma$). Grid-search was used to manipulate the values of C and $\gamma$. SVM classifier was trained and tested multiple times in order to obtain the best parameters setting that able to produce high accuracy. Lastly, the result from this survival classification is illustrated using Kaplan-Meier survival analysis curve.

## V. EXPERIMENTS AND RESULTS

This section briefly describes the experimental results obtained in four phases, data preparation phase, feature selection phase, classification phase and survival analysis phase.

### A. Data Preparation Phase

Data preparation is needed to improve the data quality by enables the data in understandable and suitable format for analysis to be performed. In this phase, data cleaning is implemented using WEKA tool to eliminate the missing and inconsistent data. The original DLBCL dataset consist 240 numbers of samples along with nine features including the class features. The outcome of this phase is the original DLBCL dataset are overwriting and leave out 220 numbers of samples.

### B. Feature Selection Phase

PSO feature selection requires parameters settings to operate. So parameters are set up before the PSO feature selection run to find the most relevant and significant features. The details on each parameters used are listed in Table I and it follows by the outcome of this phase are given in Table II.

TABLE I. PSO PARAMETERS

| Experiment | Particles (*n*) | Learning factors | |
|---|---|---|---|
| | | $C_1$ | $C_2$ |
| 1 | 5 | 0.1 | 0.1 |
| 2 | | 0.1 | 0.2 |
| 3 | | 0.2 | 0.1 |
| 4 | | 0.2 | 0.2 |
| 5 | 20 | 0.1 | 0.1 |
| 6 | | 0.1 | 0.2 |
| 7 | | 0.2 | 0.1 |
| 8 | | 0.2 | 0.2 |
| 9 | 100 | 0.1 | 0.1 |
| 10 | | 0.1 | 0.2 |
| 11 | | 0.2 | 0.1 |
| 12 | | 0.2 | 0.2 |

The setting of these parameters is based on previous studies proposed by (Rahman et al, 2009; Wang et al, 2007). Typically, we need to conduct experimentation across a range of these values to determine the best configuration of parameter setting and finally to found the optimal solution. The numbers of particles (population) were set to 5, 20 and 100. The various sizes of particles are for the comparison purposes. The maximum generation (iteration) and fitness function were set to 100 and 0.95 respectively. Meanwhile, the cognitive learning factor ($C_1$) and the social learning factor ($C_2$) were between 0.1 and 0.2 since these two values ($C_1 + C_2$) are normally limited to 4 (Rahman et al, 2009). The best solution in which each number represents one feature of DLBCL dataset at a particular iteration are listed in Table II.

The experiment does shows that the optimum features are achieved in Experiment 5, 6, 7,8,9,10,11 and 12 with features 1, 3, 4 and 5 as the best solution. The selection of best solution is finalized by compared the optimum fitness values. The combination of best solution with highest fitness value is acceptable as optimum solution. Meanwhile, Table III compared the feature length and percentage of reduction for each experiment. The feature length for experiment with optimal solution has been reduced from 8 features to 4 features with a reduction of 50%.

Table II. PSO BEST SEARCHING PROCESS ON DLBCL

| Experiment | Best Solution | Fitness Value | Feature Length |
|---|---|---|---|
| 1 | 1,9 | 0.85 | 2 |
| 2 and 4 | 1,3 | 0.85 | 2 |
| 3 | 1,3,9 | 0.85 | 3 |
| **5-12** | **1,3,4,5** | **0.8773** | **4** |

*1 represent follow-up (years), 3 represent IPI group, 4 represent germinal center B cell signature, 5 represent lymph node signature

TABLE III. RESULT BASED ON FEATURES AND
PERCENTAGE OF REDUCTION

| Experiment | Full Features | Reduced Features | Percentage Reduction (%) |
|---|---|---|---|
| 1 | 8 | 2 | 75.0 |
| 2 and 4 | 8 | 2 | 75.0 |
| 3 | 8 | 3 | 62.5 |
| 5-12 | 8 | 4 | 50.0 |

The most important thing that has been done in this phase is identify the most relevant and significant features. This experiments identify that follow-up years, IPI group, germinal center B cell signature, and lymph node signature as the most informative and significant features.

### C. Classification Phase

In order to gain insight into how PSO feature selection for classification of survival analysis works, we carried out some experiments on selected significant features that are brought from previous phase. SVM classifier has been used throughout these experiments to identify the interrelationship between selected significant features and lymphoma survival. SVM was run for 550 times of 5 folds cross validation for training and testing with respect to different parameter values of cost (C) and gamma ($\gamma$).

As suggested by (Hsu et al, 2003) we used a "grid search" on C and $\gamma$ using cross validation. It was found that SVM classifier with parameters (C $=2^7$ and $\gamma =2^{-3}$) provides higher classification accuracy with average 85.4545% correctly classified. Thus, the optimum values of parameter C and $\gamma$ are $2^7$ and $2^{-3}$ respectively. Table IV is the summarized of the performance of the DLBCL classification of survival analysis.

TABLE IV. SUMMARY DLBCL CLASSIFICATION

| Run | Survive | Dead | Correctly Classified | Incorrectly Classified | Accuracy (%) |
|---|---|---|---|---|---|
| 1 | 15 | 29 | 38 | 6 | 86.36 |
| 2 | 21 | 23 | 41 | 3 | 93.18 |
| 3 | 14 | 30 | 34 | 10 | 77.27 |
| 4 | 22 | 22 | 36 | 8 | 81.81 |
| 5 | 21 | 23 | 39 | 5 | 88.63 |
| Total | 93 | 127 | 188 | 32 | Avg= 85.45 |

### D. Survival Analysis Lymphoma Cancer Phase

As in Figure 4 and 5, we illustrated the Kaplan-Meier survival analysis curve using respective follow-up (years) feature and IPI group feature and it follows by the summary of the Kaplan-Meier result as given in Table V and VI. The curves are based on the PSO feature selection for SVM classification approach results.
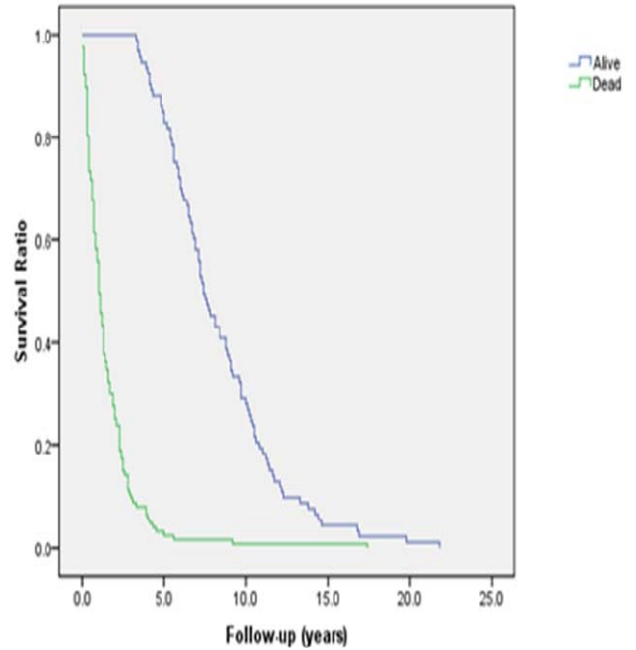


Fig. 4. Kaplan-Meier survival curve on lymphoma survival for follow-up (years) features

TABLE V. KAPLAN-MEIER RESULT ON LYMPHOMA SURVIVAL FOR FOLLOW-UP (YEARS) FEATURE

| Total N | Completed | | Censored | |
|---|---|---|---|---|
| | N of Events | Percent (%) | N of Events | Percent (%) |
| 220 | 127 | 57.7 | 93 | 42.3 |

Figure 4 presents Kaplan-Meier survival analysis curve where follow-up (which also known as survival time) is given in years. The curves for alive and dead are started with a horizontal line at a survival probability of 1.0 and then steps down to the other survival probabilities as it move from one ordered survival time to another. The surviving patients' have better survival time compared to non-surviving. The Kaplan-Meier result for follow-up feature is listed in Table V.

As present in Table V, in Kaplan-Meier survival analysis, the term *Completed* represent the patient is dead, while the term *Censored* indicated that the patient still alive. Notice that 127 out of the 220 observations are not censored representing 57.7%. From the result can conclude that 57.7% patient is defined not survived during the final follow-up. Only 42.3% patients were reported are still alive at the final follow-up visit.
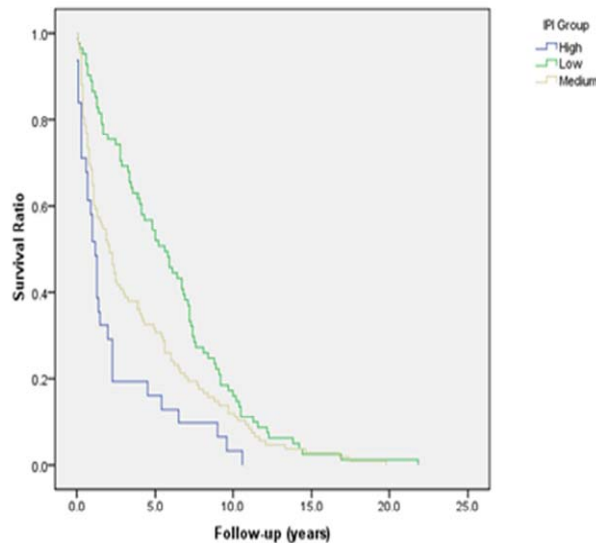
Fig.5. Kaplan-Meier survival curve on lymphoma survival for IPI group feature

TABLE VI. KAPLAN-MEIER RESULT ON LYMPHOMA SURVIVAL FOR IPI GROUP FEATURE

| IPI Group | Total N | Completed | | Censored | |
|---|---|---|---|---|---|
| | | N of Events | Percent (%) | N of Events | Percent (%) |
| Low | 81 | 27 | 33.33% | 54 | 66.66% |
| Medium | 108 | 74 | 68.51% | 34 | 31.48% |
| High | 31 | 26 | 83.87% | 5 | 16.12% |
| Overall | 220 | 127 | 57.72% | 93 | 42.27% |

Figure 5 presents the Kaplan-Meier survival analysis curve for IPI group features. The Kaplan- Meier curves were quite different with IPI group (low) having consistently better survival prognosis than IPI group (medium and high). It could be seen that IPI group (medium) having consistently better survival prognosis than IPI group (high). Note that the difference between IPI group (low) and (medium) was about the same over the time, whereas both group appeared to diverge from IPI group (high) as time increased. The Kaplan-Meier result for IPI group feature is listed in Table VI.

83.9% patient belonged to high IPI group has high probability to die during the follow-up visit, and contrary majority of the patients belonged to low IPI group are still alive which contains 66.7%. It could be observed that 57.7% patient is defined not survived during the final follow-up. Only 42.3% patients were reported are still alive at the final follow-up visit.

*E. Comparison Study*

This section, we present comparison results of PSO feature selection for SVM classification of survival analysis with survival classification using single SVM. Therefore, using the same DLBCL dataset with full data features, experiment using single SVM is carried out. As more alike in survival classification using PSO-SVM, 5-fold cross validation was applied by using the best combination input parameters ($C = 2^7$ and $\gamma = 2^{-3}$) obtained during the classification experimental in PSO-SVM previously. Single SVM without feature has least average accuracy (77.77%) compare to PSO-SVM (85.45%). It could be observed that, PSO-SVM accuracy is better than single SVM classifier. So this experiment does show the influence of number of DLBCL features over performance of classification for lymphoma survival. As stated in Table VII, PSO-SVM result for least number of features (4 features), its average accuracy is higher than classification using 8 features. These results reflect the applicability of PSO feature selection in survival classification for lymphoma patients'.

Table VII. COMPARISON RESULT FROM DIFFERENT METHODS

| | Single SVM | PSO-SVM |
|---|---|---|
| Number of features | 8 | 4 |
| Average accuracy (%) | 77.77 | 85.45 |

## VI. DISCUSSIONS

In this paper, we have proposed a PSO feature selection to overcome the limitations of classification of survival analysis in lymphoma cancer. PSO feature selection is capable of searching the optimal features for survival classification. The use of SVM classifier alone does not improve the average classification accuracy. PSO feature selection with SVM is far surpassed the efficiency of classification result. From the result, the average classification accuracy for SVM classifier with PSO feature selection performs significantly superior to the SVM classifier without feature selection. It could be seen that reducing the number of features by selecting only the significant one improved the classification accuracy. Based on the experimental result, it may appropriate to suggest feature selection for solving classification problem for survival analysis in DLBCL cancer. The performance of survival classification could be further improved by applying feature selection due to the fact that feature selection gives promising results. We had successful discovered the important role of applying feature selection that accurately classify the patient's survival.

## REFERENCES

[1]. Ahmed AA. 2005. Feature Subset Selection using Ant Colony Optimization. *International Journal of Computational Intelligence*, 53-58.

[2]. Blackwell TM. 2005. Particle Swarms and Population Diversity. *Soft Computing*, 9, 793-802.

[3]. Cortes, C. and Vapnik,V. 1995. Support Vector Networks. *Machine Learning, 20,* 273-297.

[4]. Fan C and Wan Y. 2008. An Adaptive Simple Particle Swarm Optimization Algorithm. *Control and Decision Conference,* 3067-3072.

*[5].* Fox J. 2006. Introduction to Survival Analysis. *Lecture Notes.*

[6]. Geetha K, Thanushkodi K and Kumar AK. 2008. New Particle Swarm Optimization for Feature Selection and Classification of Microclacifications in Mammograms. *International Conference on Signal Processing,Communication and Networking,Madras,* 458-463.

[7]. Hamdan H and Garibaldi JM. 2012. A Framework For Automatic Modeling of Survival Using Fuzzy Inference. *IEEE International Conference on Fuzzy Systems*, 1-8.

[8]. Hasan H and Tahir NM. 2010. Feature selection of Breast Cancer Based on Principal Component Analysis. *International Colloqium on Signal Processing and Its Application.*

[9]. Hsu CW, Chang CC and Lin CJ. 2003. A Practical Guide To Support Vector Classification.

[10]. Jamian JJ, Mustafa MW. Mokhlis H and Abdullah MN. 2012. Comparative Study on Distributed Generator Sizing Using Three Types of Particle Swarm Optimization. *International Conference on Intelligent Systems, Modelling and Simulation*, 131-136.

[11]. Jensen R. 2005. Combining Rough and Fuzzy Sets For Feature Selection. *Phd Thesis*, University of Edinburgh.

[12]. Kennedy J and Eberhart RC. 2001. *Particle swarm intelligence* Morgan Kaufmann Publishers, San Francisco, USA.

[13]. Kennedy J and Spears WM. 1998. Matching Algorithms to Problems: An Experimental Test of the Particle Swarm and Some Genetic Algorithms on the Multimodal Problem Generator. *IEEE World Congress on Computational Intelligence*, 78-83.

[14]. Kumar AM and Gopal M. 2010. A hybrid SVM Based Decision Tree. *Pattern Recognition,*43, 3977-3987.

[15]. Liu Y, Wang G, Chen H,Dong H, Zhu X and Wang S.2011. An Improved Particle Swarm Optimization for Feature Selection. *Journal of Boinic Engineering*, 8.

[16]. Liu Y. 2010. Survival Analysis of Breast Cancer. *Master Thesis*, University of Victoria.

[17]. Mishra D and Sahu B. 2011. Feature selection for Cancer Classification: A Signal-To-Noise Ratio Approach. *International Journal of Scientific and Engineering Research*, 2.

[18]. Omar N, Jusoh F, Ibrahim R and Othman MS. 2013. Review of Feature Selection for Solving Classification Problems. *Journal of Information System Research and Innovation*. 3.

[19]. Othman MS. 2008. An Automatic Web Information Resources Classification Using Extraction and Machine Learning Approach, *Phd Thesis*, Universiti Kebangsaan Malaysia.

[20]. Rahman, SA, Bakar AA and Hussein ZAM. 2009. Filter-wrapper approach to feature selection using RST-DPSO for mining protein function. *Conference on Data Mining and Optimization*, 71-78.

[21]. Shahamatnia E and Ebadzadeh MM. 2011. Application of particle swarm optimization and snake model hybrid on medical imaging. *IEEE Third International Workshop On Computational Intelligence In Medical Imaging*, 1-8.

[22]. Sharkawy RM,Ibrahim K, Salama MMA and Bartnikas R. 2011. Particle swarm optimization feature selection for the classification of conductiong particles in transformer oil. *IEEE Transaction on Dielectrics and Electrical Insulation*, 18.

[23]. Tu CJ, Chuang LY, Chang JY and Yang CH. 2007. Feature selection using PSO-SVM. *IAENG International Journal of Computer Science*, 33.

[24]. Wang X, Yang J, Teng X, Xia W and Jensen R. 2007. Feature Selection Based on Rough Sets and Particle Swarm Optimization. *Pattern Recognition Letters*, 28.

[25]. Wei Z, Ling H.Y, Tao ZH and Jing TW. 2008. Enhancing the Particle Swarm Optimization Based on Equilibrium of Distribution. *Control and Decision Conference,* 285-289.

*[26].* Yang CS, Chuang LY. Li JC, and Yang CH. 2008. Chaotic Maps in Binary Particle Swarm Optimization for Feature Selection. *IEEE Conference on Soft Computing in Industrial Application.*

[27]. Zhanshan ZM and Survival AWK. 2008. Survival Analysis Approach to Reliability, Survivability and Prognostics and Health Management. *IEEE Aerospace Conference, 1-20.*

[28]. ZhongXin, D. 2011. The Application of Support Vector Machine in Survival Analysis. *International Conference on Artificial Intelligence, Management Science and Electronic Commerce*, 6816-6819.