# An Adaptive Binary PSO to Learn Bayesian Classifie for Prognostic Modeling of Metabolic Syndrome

Satchidananda Dehuri
Dept. of Information and
Communication Technology
Fakir Mohan University
Vyasa Vihar,
Balasore - 756019
Orissa, India
satchi.lapa@gmail.com

Rahul Roy
School of Computer
Engineering
KIIT University
Bhubaneswar
Orissa-751024, India
link2rahulroy@gmail.com

Sung -Bae Cho
Soft Computing Laboratory
Dept. of Computer Science
Yonsei University
262 Seongsanno,
Seodaemun-gu
Seoul 120-749, South Korea
sbcho@cs.yonsei.ac.kr

## ABSTRACT

The metabolic syndrome is a combination of medical disorders that have become a significant problem in Asian countries due to the change in lifestyle and food habits. Thus a prognostic model can help the medical experts in diagnosis of the disease. Learnable Bayesian classifier by Adaptive Binary Particle Swarm Optimization (ABPSO) provides a robust formalism for probabilistic modeling that can be used as a predictive tool in medical domain. In this paper, we adopt an ABPSO for adapting the weights of the learnable Bayesian classifier that provides a maximum prediction accuracy and can exhibit an improved capability of removing spurious or little important attributes and help the medical experts in identifying the basis for the disease. Experiments have been conducted with the dataset obtained in Yonchon Country of Korea, and the proposed model provides better performance than the other models.

## Categories and Subject Descriptors

H.2.8 [**Database management**]: Database Applications—
*Data mining*

## General Terms

Theory Algorithms, Performance, Design, Experimentation.

## Keywords

Metabolic syndrome, Naive Bayesian classifier, Feature selection, k-means, Binary particle swarm optimization

## 1. INTRODUCTION

The metabolic syndrome, characterized by a constellation of metabolic disorders including dyslipidemia, elevated blood glucose, hypertension, and obesity, is recognized as a major looming epidemic of the $21^{st}$ century. The correlation between cardiovascular risk factor and metabolic syndrome is reported in studies [6]. As reported in [6, 7, 8], the disorder is more common in adults of 20 years and above. Approximately 34% of adults met the criteria for metabolic syndrome. Males and females $40 \sim 59$ years of age were about three times as likely as those $20 \sim 39$ years of age to meet the criteria for metabolic syndrome. In Asian countries, it has become a significant problem lately due to the change in dietary habit and life style. Approximately $5 \sim 17\%$ of the males and females of the age $20 \sim 72$ years have the prevalence of metabolic syndrome diseases as reported in [7]. Due to the advent of metabolic syndrome as a epidemic, it has become an interest of study from many of the researchers around the world [5].

Bayesian classifier [1] is a statistical method for solving the problem of classification. It has emerged as a simple,yet effective classifier and has been widely applied to complex domains. This classifier has been extensively used as a predictive model in their medical domains and have shown a high performance. The successful application of the Bayesian classifier to the diagnosis and prognosis of diverse diseases have motivated us for adopting this technique for designing a prognostic model for the prediction of metabolic syndrome.

This paper deals with a problem that predicts the metabolic syndrome with the dataset obtained in Yonchon County of Korea. This paper presents a statistical analysis of the dataset and makes a prognostic model using a learnable Bayesian classifier where the learning vectors are determined using an ABPSO algorithm, which is a binary version reported in [4]. Using this real world dataset, the performance of some of the state-of-the-art predictive model along with ABPSO aided learnable Bayesian classifier is studied in this paper. Also the impact of the parameters used in ABPSO on the prediction accuracy is studied in this paper.

The rest of the section is organized as follows: Section 2 provides a brief review on the extensive use of Bayesian classifier in the medical domain. Section 3 provides a basics of naive Bayesian classifier and its extension, and an adaptive Binary PSO aided learnable Bayesian classifier. Section 4 describes the metabolic syndrome problem and Section 5 provides the predictive model for the classifier. Section

6 provides a statistical study of the dataset and analytical results of the evaluation of the classifier on the dataset. Finally conclusion and future research direction are presented in Section 7

## 2. RELATED WORKS

Bayesian classifier has been robustly used in the medical domain for the prognostic modeling of various diseases. Here, a short review on some of the application of the Bayesian classifier for prediction in medical domain has been discussed. In 2001 Raymer et al., [10] proposed a hybrid algorithm based on the Bayes discriminant and evolutionary algorithm. They employed their algorithm for feature selection and extraction to isolate the salient features from large medical and biological dataset. Blanco et al. [2] in 2005, induced filters and wrappers based on feature subset set selection, in Bayesian classifiers for selection for predictive modeling of survival of cirrhotic patients treated with TIPS. In 2006, Cruz-Mesia and Quintana [3] proposed a Bayesian classifier based on the hierarchical model for non-linear longitudes profiles. They applied the classifier for predicting the pregnancy outcomes from longitudinal $\beta$-hCG profiles.

In 2007, Wiggins et al. [12], proposed a genetic algorithm based Bayesian classifier for predicting ECG-based age classification. They applied genetic algorithm for finding the structure of the Bayesian network that determines the conditional probability among the features. In the same year, Abraham et al. [9], made a comparative study of the statistical classifier on 21 medical dataset and have shown that the performance of naive bayes classifier (MDL discretized) is better compared to other classifiers. In 2008, Taieb et al. [11], proposed a iterative Bayesian approach for tumor analysis in liver. They provide a method and validation study for the nearly automatic segmentation for the liver tumors.

## 3. PRELIMINARIES

### 3.1 Learnable Bayesian Classifie

The naive Bayesian classifier is a probabilistic approach to classification. Given an unclassified object $X = (x_1, x_2..., x_n)$ the classifier predicts that $X$ belongs to the category having the highest posterior probability conditioned on $X$. Specifically, this classifies object $X$into category $C_i$ if and only if

$$P(C_i|X) > P(C_j|X) \quad \forall i \neq j$$

The $P(C_i|X)$ is calculated using the bayes theorem.

$$P(C_i|X) = \frac{P(X|C_i)P(C_i)}{\sum_i P(X|C_i)P(C_i)}$$

The training sample is used to determining the required probabilities. The $P(C_i)$ is calculated as $P(C_i) = \frac{q_i}{q}$ where q is the total number of samples and $q_i$ is number of samples of class$C_i$.

The crucial part of the formulation of the Bayesian model is the determination of $P(X|C_i)$. The determination of $P(X|C_i)$ can be extremely computationally expensive and may require large training samples. The naive Bayesian model makes the simplifying assumption that

$$P(X|C_i) = \prod_{k=1}^{n} P(x_k|C_i)$$

This assumption, called attribute conditional independence, greatly simplifies the calculation.

However, naive Bayesian classifier does not provide any assurance for accurate classification of training dataset. The training set is only used for determining the probabilities. Thus Yager et al. [13], provided an extension for naive Bayesian classifier by introducing learning weights in the model. The formulation for the $P(X|C_i)$ is shown in equation 1

$$P(\overrightarrow{x}|C_i) = \sum_{j=1}^{n} w_j \left( \prod_{k=1}^{j} P(x_k|C_i) \right) \quad (1)$$

Where $P(x_k|C_i)$ is the $k^{th}$ largest of the $P(x_j|C_i)$ and $w_j = [0,1]$ and $\sum w_j = 1$

The introduction of this weights provides with the additional degree of freedom in the form of the associated weights. However, these weights are very sensitive to local optima. So, we need a good heuristic method to find these weights. PSO being an population based heuristic method has a less chance of getting trapped in local optima and is a good alternative for finding these weights.

### 3.2 Adaptive Binary Particle Swarm Optimization

The adaptive binary particle swarm optimization [4] is a combination of the best effort of both the adaptive PSO and binary PSO to explore the continuous and the discrete search space simultaneously. In this approach, each particle of the swarm is initialized using values from both continuous and discrete domains and search for optima by updated in each iteration. The search of the particle is guided by using two best values. The first one is the local best position of the particle and the second one is the global best position of the swarm.

The inertia weight $(w)$ is the most important parameter that moves the particle towards the optimal position. Thus, to increase the search ability the particles flight should be controlled by the objective functions. The particle which is closer to the optimal point should move slowly as compared to the other particle. This movement of the particle can be controlled using different $w$ values according to their rank between $w_{min}$ and $w_{max}$ as give in following equation 2:

$$w_i = w_{\min} + \frac{w_{\max} - w_{\min}}{T_{pop}} \times rank_i, \quad (2)$$

where $T_{pop}$ is the total size of the swarm. The velocity of the particle,both for continuous and discrete space, is updated using the equation 3.

$$v_{id}(t+1) = wv_{id}(t) + c_1 rand_1()(p_{id}(t) - x_{id}(t)) + \\ c_2 rand_2()(p_{gd}(t) - x_{id}(t)), \quad (3)$$

However, the position updation of the particle represented by continuous vector is done using equation 4 and discrete vector is updated using equation 5 respectively.

$$x_{id}(t+1) = x_{id}(t) + v_{id}(t+1) \quad (4)$$

$$S(v_{id}) = 1/(1 + exp(-v_{id}))$$
$$x_{id} = \begin{cases} 1 & \text{if} \quad S(v_{id}) > rand() \\ 0 & \text{otherwise} \end{cases} \qquad (5)$$

## 4. METABOLIC SYNDROME PROBLEM

The definition of the metabolic syndrome was provided by the National Cholesterol and Education Program, Adult Treatment Panel III (ATP III). It requires the presence of three or more of the following components [6].

- Abdominal obesity (waist circumference $> 102$ cm in men and $> 88$ cm in women),

- Hypertriglyceridemia ( $= 150$ mg/dL),

- Low high density lipoproteion (HDL) cholesterol$< 40$ mg/dL in men and $< 50$ mg/dL in women),

- High blood pressure (systolic $= 130$ mmHg or diastolic $= 80$ mmHg) and

- High fasting glucose ($> 110$ mg/dL).

Since this original standard is not appropriate for Asian, we have used modified the definition for Asian of the abdominal obesity (waist circumference$> 90$ cm in men and $> 80$ cm in women) in this paper [5].

## 5. PROGNOSTIC MODELING OF METABOLIC SYNDROME

Figure 1 represents the illustrative pictorial diagram of the the prognostic model. This process is divided into three parts. The first part involves pre-processing of the dataset where the medical domain knowledge is applied for better design of the model. The second phase involves the construction of the learnable naive Bayesian classifier with their weight adjustment using the training set. Finally the process ends in third phase with prediction of the class label for the unknown samples.

### 5.1 Preprocessing

During the pre-processing stage,we have analyzed the dataset in order to provide a summarization of the dataset and study the relationship among the attributes. The dataset is a mixed dataset with two discrete attribute and sixteen continuous attributes.

The dataset is also tested for the presence of outlier using the k-means based outlier detection technique. The k-means algorithm is applied for clustering the dataset where k is the number of clusters specified in the dataset. Then the Euclidean distance of the centroid from the every point of the cluster is calculated. Then the points with distance which are greater than the threshold value are declared as outlier points. The threshold value is given as the 1.5 times of the average of the distance.

### 5.2 ABPSO Learnable Bayesian classifie

During the construction of classifier, we find the mean and standard deviation of the continuous attribute belonging to each class and the probability $P(x_j|C_i)$ for the discrete attributes using the training set. Then the optimal learning vector used in equation (1) is found using the adaptive binary PSO.
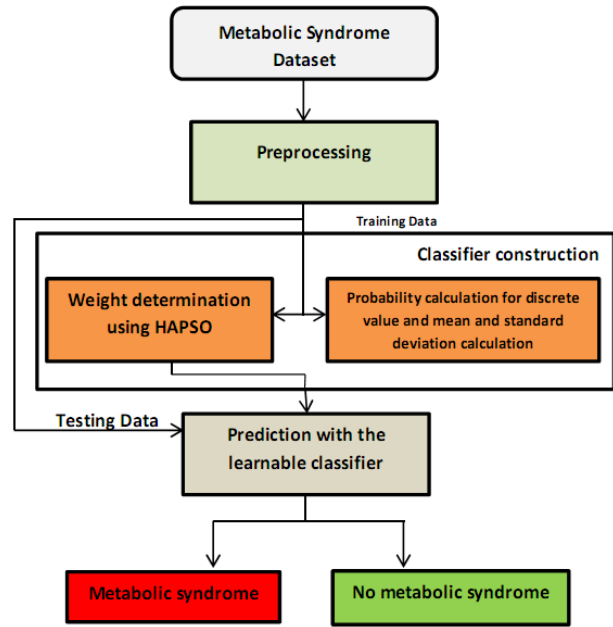


Figure 1: Pictorial representation of the prognostic model

The particle representation is shown in Figure 2. The particle is represented as $2 \times n$ matrix, where n is the number of attributes. the first vector of the matrix represents the weights $w_i$ and the second vector represents the mask bit $m_i$ i.e., the inclusion or exclusion of the weights in the learning process.



Figure 2: Particle Representation

The weights used for updating the particle's position is calculated using equation (2). The velocity and the position of the particle is updated using equations (3-5). The particle with the highest fitness value is selected as the guide for the search process ($gbest$).

During the execution of the algorithm, each individual is passed through the classifier for evaluation and the cost is computed based on the classification accuracy obtained from the parameterized generalized bayesian formulation in classifying the known set of samples of known class. The HAPSO algorithm seeks to maximize the cost function. The mathematical formulation of the cost function is given in equation (6):

$$f(\overrightarrow{x}) = A_c \times CL_{accu} + \frac{A_f}{Sfr}, \qquad (6)$$

where $CL_{accu}$ is the classification accuracy, $A_c$ is the weight factor associated with the classification accuracy, $A_f$ is the weight factor associated with the selected number of weights and $Sfr = \sum m_i \times rank(w_i)$. The weights vectors are as-

signed rank based on the rank of $P(x_j|C_j)$. The value for the $A_c$ and $A_f$ are empirically found.

The designed classifier efficiency is evaluated based on the predictive accuracy of the test dataset. The algorithm for the learnable bayesian classifier is shown in Algorithm 1.

---

**Algorithm 1** Pseudo code for learnable Bayesian classifier

calculate the mean and standard deviation for continuous values and prior probability for discrete values from the training set and save it in a file.
Initialize the swarm $S$ with swarmsize
**for** each particle of the $S$ **do**
    Initialize the velocity $v_i$
    Initialize the *pbest* $p_i$
**end for**
Evaluate the fitness $F = fitness(S)$;
Select the *gbest* as the particle from swarm $S$ with maximum fitness value $F_i$
**while** stopping criterion is not met **do**
    **for** each particle **do**
        Find the weight $w$ using equation (2)
        Update the velocity using equation (3)
        **if** particle is real valued vector **then**
            Update particle using equation(4)
        **else**
            Update particle using equation(5)
        **end if**
    **end for**
    Evaluate the fitness $F_{new} = fitness(S)$;
    **for** each particle i **do**
        **if** $F_{new}(i) > F(i)$ **then**
            update the *pbest* $p_i$
        **end if**
    **end for**
    **if** $max(F_{new}) > max(F)$ **then**
        update the gbest with the particle having maximum value of $F_{new}$
    **end if**
    $F = F_{new}$
**end while**

---

**function Fitness(S)**

    **for** each particle i **do**
        find the classification accuracy and the rank
        $[accuracy, rank] = $classification(training set, $particle_i$)
        find the fitness based on accuracy and rank using equation(6)
    **end for**
    return fitness

---

# 6. EXPERIMENTAL STUDY

In this section, a detailed analysis of the data set is provided followed by the performance metrics used for evaluating the classifiers, experimental results, discussion and analysis of the results on the dataset.

## 6.1 Analysis of Dataset

The metabolic syndrome dataset has 1135 data tuples with 18 attributes and two class levels. Class 0 and class 1 represents the not occurrence and occurrence of the metabolic

---

**Function classification(Dataset, learning vector)**

    load file containing the prior probability, mean and standard deviation.
    **for** each sample of the dataset **do**
        find the probability $P(x_j|C_i)$
        sort $P(x_j|C_i)$ in descending order and find the rank.
        Calculate the $P(X|C_i)$using the equation (1)
        predict the class label of the sample as $C_i$ which has the maximum $P(X|C_i)$.
    **end for**
    find the true positive classification
    $accuracy = \frac{truepositive}{|dataset|}$
    return [accuracy,rank]

---

**Table 1: Statistical Analysis of the Dataset**

| Attribute | Mean | Median | Variance | Standard Deviation |
|---|---|---|---|---|
| Age | 53.78238 | 54 | 141.9394 | 11.91383 |
| Sex | 1.5691 | 2 | 1.1383 | 1.0669 |
| fasting glucose | 101.1507 | 100 | 102.2495 | 10.1185 |
| Post Prandial Glucose | 107.4907 | 104 | 784.2887 | 28.00515 |
| Height | 157.9436 | 158 | 73.1933 | 8.555308 |
| Weight | 60.41498 | 60 | 99.30004 | 9.964941 |
| AC | 81.53304 | 82 | 143.6427 | 11.9851 |
| HC | 93.6326 | 94 | 120.6007 | 10.98183 |
| HT | 1.882 | 2 | 1.1837 | 1.0879 |
| GOT | 21.11 | 18 | 317.2628 | 17.81187 |
| GPT | 18.31718 | 14 | 298.9496 | 17.81187 |
| Cholesterol | 158.5894 | 155 | 959.2746 | 30.9216 |
| Triglyceride | 149.3859 | 122 | 11193.06 | 105.7972 |
| HDL cholesterol | 37.51982 | 36 | 142.7712 | 11.94869 |
| Body mass index | 24.1815 | 23.9 | 10.87968 | 3.298435 |
| weight-height ratio | 0.870088 | 0.87 | 0.004692 | 0.068497 |
| Systolic blood pressure | 125.3604 | 120 | 472.3239 | 21.73301 |
| Diastolic blood pressure | 81.33833 | 80 | 193.5631 | 13.91269 |

syndrome disease respectively. Each class have 612 and 523 data tuples respectively. The measures for central tendency and dispersion of each attribute data are shown in Table 1.

The deviation around mean is maximum in Triglyceride attribute and for all other attribute the dispersion does not exceed the more than 31. The box plot of the each have shown that attribute Triglyceride has an highest maximum value which greater than $1.5\times$ Inter Quartile Range (IQR) of the attribute. This has led to its abnormal rise in deviation. The Boxplot for attributes identifying each classes showed that class 0 has less outlier points as compared to class 1. Some boxplots depicting the fact is shown in Figure 3. Also attributes identifying class 0 has better spread as compared to class 1. The systolic blood pressure attribute has a single value (i.e., the median value) for identifying the class 1.

Pearson correlation coefficient is also calculated to evaluate the bonding among the attributes. We could obtain a height positive correlation of 0.8530 among GOT and GPT and a highest negative correlation of -0.7020 with the at-

tribute sex and height. The Scatter plot for both the positive correlation and negative correlation is shown in Figure 6.1. Further, analysis of each attribute's correlation with other attributes is summarized below. Age attribute have little or no correlation with any of the other attributes. Its highest positive correlation obtained is 0.2051, with attribute fasting blood sugar and highest negative correlation obtained is -0.2586 with attribute height. The attribute AC and HC, weight, body mass index and Systolic blood pressure and diastolic blood pressure shows strong positive correlation with correlation coefficient in the range of [0.75, 0.8489]. Except the sex and height, we do not find much higher negative correlation among the other attributes.

This dataset is good for classification problem but we need to take some measures for the outlier points especially for attribute Triglyceride and AC. There is not much dependency in the dataset which is evident from the correlation analysis except for few of the attributes.

## 6.2 Performance Metrics

The performance metrics used for evaluating the performance of the classifier for predicting the metabolic syndrome are give below:

*Precision:*.
Precision measures the percentage of tuples that are correctly classified. The measure is defined as

$$Precision = \frac{TP}{TP + FP}$$

where $TP$ is the number of true positives and $FP$ is the number of false positives(tuples of the negative class incorrectly classified as positive class).

*Recall:*.
Recall is the fraction of correct instances among all instances that actually belong to the relevant subset. Mathematically, it is defined as

$$Recall = \frac{TP}{TP + FN}$$

where $FN$ is the false negative (i.e., tuples belonging to positive class is classified as negative class).

*$F_1$-score:*.
The $F_1$ score (also F-score or F-measure) is a measure of a test's accuracy. It is a harmonic mean of precision and recall and is defined as

$$F_1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

## 6.3 Experimental Setup

The experiment is carried out using Matlab 2009 in multicore system with core2duo processor, 2 GB ram under windows OS. The parameter setting used for the implementation are shown below:
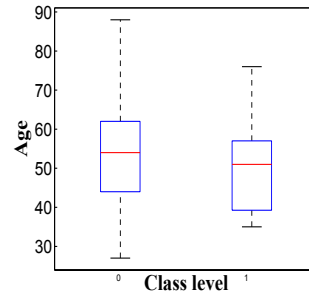Population Size: 50
Weight $w_i$ for swarm movement: 0.4-0.9
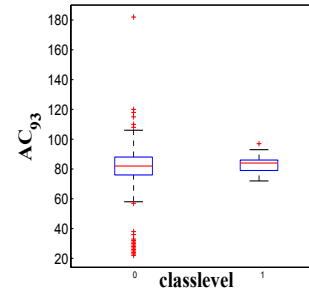Personal Influence Factor $c_1$: 1.8
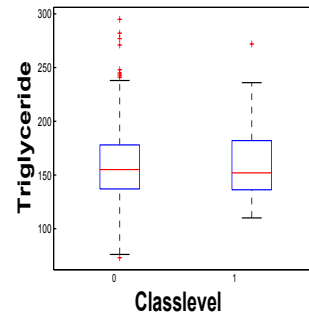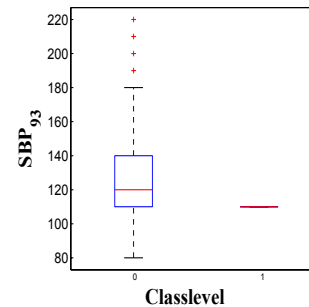Global Influence Factor $c_2$: 2.5
$A_f$ :10
$A_c$ : 20



(a)



(b)



(c)



(d)

Figure 3: Boxplot for the attributes age, AC,triglyceride , and systolic blood pressure for each class level.
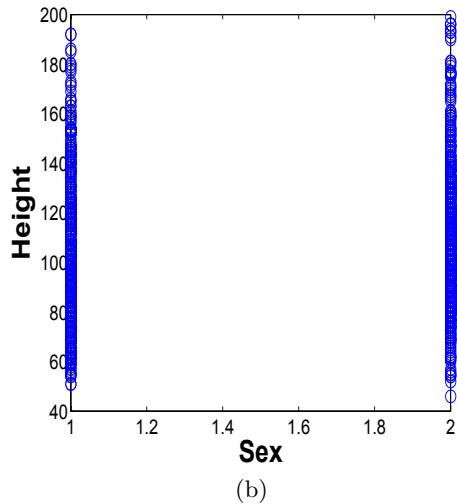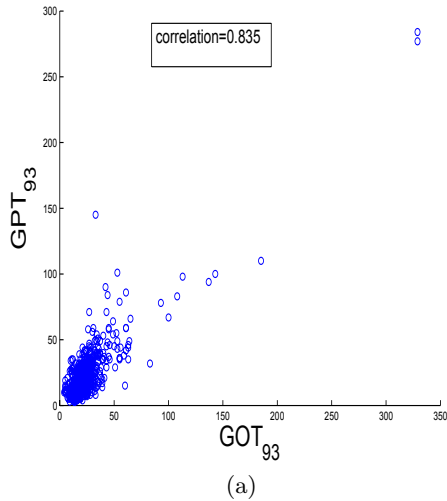
499

|                 | Accuracy (%) | Attributes |
|-----------------|--------------|------------|
| With Outlier    | 72           | 11         |
| without Outlier | 76.14        | 11         |

Number of Iteration: 80

The dataset is divided into 10 disjoint for 10 cross fold validation. For each fold one set is used as testing set and remaining set is used for training/tuning the classifier.

## 6.4    Results and Analysis

First the experiment is conducted using the dataset with outlier points and after removing the outlier points the experiment has been repeated again. The results are shown in the Table 2. By removing the outlier points, we obtained a significant improvement predictive accuracy. Also, the classifier performance is compared with some of the other classifiers based on precision, recall and $F_1$-score. To calculate each of the metric, a confusion matrix for each of the classifier is determined and then the performance metrics is calculated. The results are given in Table 3. The precision and recall of the ABPSO based learnable classifier is higher compared to other classifier. The $F_1$-score clearly shows that the test accuracy of ABPSO based learnable Bayesian classifier is better compared to other classifier.

**Table 3: Performance evaluation of Other classifier on metabolic syndrome dataset**

| Algorithm | Accuracy (%) | Precision (%) | Recall (%) | $F_1$-score |
|-----------|--------------|---------------|------------|-------------|
| ABPSO learnable bayesian classifier | 76.14 | 77.215 | 85.165 | 0.81 |
| C45 | 68.26 | 63.40 | 61.98 | 0.6268 |
| ID3 | 68.72 | 67.03 | 69.73 | 0.6835 |
| K-NN | 71.28 | 71.15 | 72.85 | 0.7198 |
| Multilayer perceptron | 72.94 | 75.13 | 73.68 | .7439 |
| SVM | 75.96 | 77.13 | 75.18 | 0.7614 |

## 7.    CONCLUSION

This paper present an ABPSO learnable Bayesian classifier for predicting the metabolic syndrome disease. A detailed analysis of the data set showed that few attributes have outlier points which can be a cause for misclassification. It is shown in [8] that optimal accuracy is 72%. This is found to be true with outlier points even in ABPSO. However, when the outlier detection technique was applied, we found that there is a increase in accuracy. Thus the incorporation of the clustering based outlier detection proved to be a contributing factor in the increase of accuracy. However, this performance evaluation is limited to single dataset. we need to evaluate the performance with other dataset.



**Figure 4: Correlation analysis of GPT, GOT and Sex and height**

## 8. REFERENCES

[1] M. Bayes and M. Price. An Essay towards Solving a Problem in the Doctrine of Chances. By the Late Rev. Mr. Bayes, F. R. S. Communicated by Mr. Price, in a Letter to John Canton, A. M. F. R. S. *Philosophical Transactions*, 53:370–418, 1763.

[2] R. Blanco, I. Inza, M. Merino, J. Quiroga, and P. Larrañaga. Feature selection in bayesian classifiers for the prognosis of survival of cirrhotic patients treated with tips. *Journal of Biomedical Informatics*, 38(5):376 – 388, 2005.

[3] R. de la Cruz-Mesía and F. A. Quintana. A model-based approach to bayesian classification with applications to predicting pregnancy outcomes from longitudinal $\beta$-hcg profiles. *Biostatistics*, 8(2):228–238, May 2006.

[4] S. Dehuri, B. K. Nanda, and S.-B. Cho. A hybrid apso-aided learnable bayesian classifier. In *IICAI*, pages 695–706, 2009.

[5] M. K. Moon, Y. M. Cho, K. S. Lim, P. S., and H. K. Lee. Metabolic syndrome. *The Korean Society of Endocrinology*, 18:105–116, 2003.

[6] L. Mykkänen, J. Kuusisto, K. Pyörälä, and M. Laakso. Cardiovascular disease risk factors as predictors of type 2 (non-insulin-dependent) diabetes mellitus in elderly subjects. *Diabetologia*, 36:553–559, 1993.

[7] P. Nestel, R. Lyu, L. P. Low, W. H.-H. Sheu, W. Nitiyanant, I. Saito, and C. E. Tan. Metabolic syndrome: recent prevalence in east and southeast asian populations. *Asia Pac J Clin Nutr*, 16(2):362–327, 2007.

[8] H.-S. Park and S.-B. Cho. An efficient attribute ordering optimization in bayesian networks for prognostic modeling of the metabolic syndrome. In *ICIC (3)*, pages 381–391, 2006.

[9] J. Ranjit Abraham, Simha and S. Iyengar. Using probabilistic classifiers for medical data mining. white paper, department of Computer Science, Louisiana State University, Baton Rouge, USA, 2006.

[10] M. L. Raymer, L. A. Kuhn, and W. F. Punch. Knowledge discovery in biological datasets using a hybrid bayes classifier/evolutionary algorithm. In *Proceedings of the 2nd IEEE International Symposium on Bioinformatics and Bioengineering*, pages 236–245, Washington, DC, USA, 2001. IEEE Computer Society.

[11] Y. Taieb, O. Eliassaf, M. Freiman, L. Joskowicz, and J. Sosna. An iterative bayesian approach for liver analysis: tumors validation study. In *In Proc. of the MICCAI 3D Segmentation in the Clinic: A Grand Challenge II workshop, in MICCAI'08*, 2008.

[12] M. Wiggins, A. Saad, B. Litt, and G. Vachtsevanos. Evolving a bayesian classifier for ecg-based age classification in medical applications. *Applied Soft Computing*, 8(1):599 – 608, 2008.

[13] R. R. Yager. An extension of the naive bayesian classifier. *Information Sciences*, 176(5):577 – 588, 2006.