

# The Classification Techniques on Medical Data to Predict Heart Disease

Gayam Suhasini  
*M.Tech Student*

*Department of Information Technology  
VR Siddhartha Engineering College, Vijayawada, India*

Vemuri Sindhura  
*Assistant Professor*

*Department of Information Technology  
VR Siddhartha Engineering College, Vijayawada, India*

Sandeep Y  
*Assistant Professor*

*Department of Information Technology  
VR Siddhartha Engineering College, Vijayawada, India*

**Abstract:** Analysis of therapeutic information is very test in context of incremental development of properties and different parameters. Once the information is gathering there is specific farthest point to wind up in getting the information as tuples or ascertaining the recurrence of the combinational credits concerning age, ailment, sexual orientation. The most compelling motivation is to push the mysterious maladies in examination. In substantial information related information spotting of the fancied infection information is excessively muddled. So we utilize KNN with Euclidean separation component and choice tree with ordinary and upgraded model concerning given characteristics. In KNN Euclidean separation produced for all tuples and rank will be accommodate new characterization of new set and result created in light of k worth as closest neighbors. Be that as it may, for examination is between both above said as for investigation of time multifaceted nature for order. The principle characterization is done on tremendous and element patch information. So the fundamental arrangement is done on KNN methodology and immediate arrangement trees make up by utilizing NAE approach (typical and improved choice tree structures).

**(Key words:** KNN, NAE, decision tree, Classification, Euclidean distance)

## I. INTRODUCTION

The fundamental point of this work is to accomplish the achievable investigation on gigantic restorative information with deference every one of the traits in the information sets. The work is on various sorts of properties like cholesterol and BP. In KNN characterization therapeutic properties are ordered into various classes. Every one of these classes fall under various combinational and restrictive extents as for restorative entanglements. In kNN finding the closest neighbors concerning new classification of qualities and arranges and finding the outcome as classes. This outcome is relies on upon estimation of K and neighbors of new tuples Euclidian separation. In choice tree we have considered typical and improved choice tree. The principle contrast is to diminish the tree progressive system when contrast with improved tree with ordinary tree. In Normal choice tree and upgraded choice tree we consider same ascribes to assemble tree to anticipate discover the danger variable of coronary illness. Upgraded choice tree produces precise result contrasted with Normal choice Tree. In KNN characterization taking into account quality extents for presentation however for this situation displaying wanted result and time many-sided quality of order is completely taken a toll utilization. So our work is relocated in receptive way as choice tree in low and abnormal state of showing the sought yield. [1][3]

## II. PROPOSED METHOD

The proposed framework as appeared in fig1 comprises of 3 phases i.e., crude information is gathered from lab data system, Diagnosis, Procedures, Pharmacy, Procedure Notes, Nurse Records. The crude information is separated and obliged credits to perform order are recovered. Machine learning calculations are connected on medicinal information gathered to anticipate heart disease..kNN calculation, Decision tree, and Enhanced Decision tree are the procedures used to characterize the restorative information. In kNN calculation classifiers are produced after that Euclidean separation is computed to discover the class of the new info therapeutic records. Get the k worth to create number of classes identified with new records. In choice Tree perusing dataset is performed to recover as far as possible to manufacture root hub and cholesterol and overweight reaches are considered to fabricate further level hubs. The check of the considerable number of records which have high hazard element of coronary illness are shown with Aadhaar Id's of the patients. In Enhanced Decision Tree perusing of tuples is performed and sexual orientation is taken as root hub and the pair of ascribes are considered to work next level hubs like overweight, cholesterol and mid-section torment for female patients and smoking ,liquor and cholesterol are considered as leaf credits to anticipate hazard element of coronary illness. Last tally of all the patients who have high hazard element are shown as for their Aadhaar numbers. The best strategy is discover of the three machine learning calculation in view of the exactness estimations.

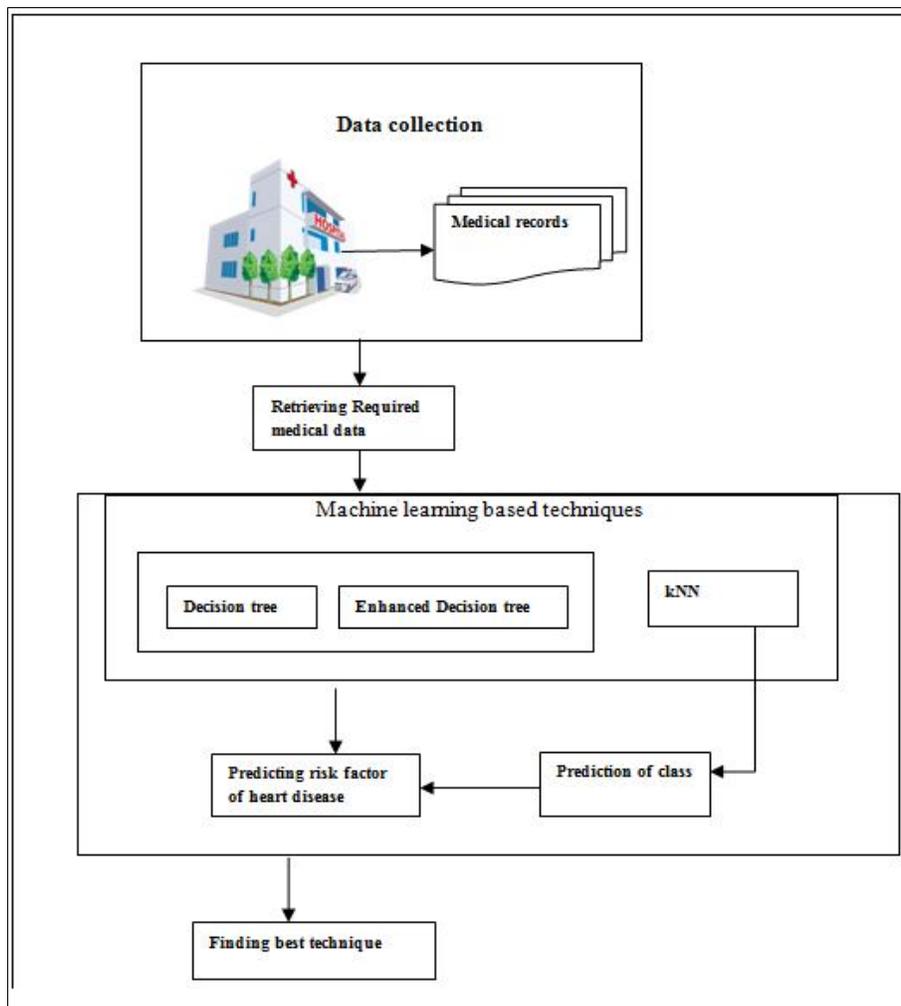


Fig 1: Architecture of proposed system

## III. MACHINE LEARNING TECHNIQUES

Mining is utilized to make insightfulness and utilization of data. Information revelation in Data is non-minor procedure of ordering usable, one of a kind, hypothetically significant and understandable methodologies in information. Information mining comprises of more than social event of taking care of information. It likewise contains investigation and expectation. Information mining comprises a few calculations, for example, Decision

Tree, kNN, Naive Bayesian, Support Vector Machines, Apriori calculation. Every one of these calculations will perform well in restorative conclusion, social insurance framework and drug. In this paper Decision Tree, Enhanced Decision Tree and kNN calculations are utilized to perform examination on therapeutic information to anticipate class objects.

### 3.1 kNN

Medicinal information is so progressively developed source as datasets of data from different doctor's facilities as tuples as patient records. At the point when information sets mined, the shrouded data in these datasets is a major asset group for medicinal examination and presentation. Every one of the information contains diverse examples and close arrangement of relations, which can come about for better analysis. The fundamental extreme undertaking in finding and characterization of these examples and relations regularly uncovered in IT. Research has been moved in therapeutic determination to discover heart maladies, lung ailments and different thyroid issues in light of the information gathered on influenced patients. In any case, there is some confinement to particular space frameworks that can research illnesses limited to their applicable operations. By and large the bore of k-closest neighbourhood (k-NN) classifier is thoroughly rely on upon the separation metric used to perceive k closest neighbors of the question focuses. The standard Euclidean separation is usually utilized as a part of general information analysis.[10]The research utilizes enormous capacity of data with the goal that finding is completely in view of verifiable information can be made. The objective is on figuring the likelihood of produced of a specific affliction by utilizing another calculation. This kNN calculation immeasurably builds the exactness of finding. This methodology can be utilized to redesign the mechanized in analysis., which blends determination of numerous infections indicating comparative characteristics and symptoms.[2]

$$Distance = \sqrt{(an - ap)^2 + (cn - cp)^2} \quad (1)$$

**Algorithm:**

```

EUCLIDIAN DISTANCE CALCULATION AND RANKING
Cp = GETn()
Bp = GETn()
Σ ED = 0
for each n in C
ED(n) =  $\sqrt{(C_n - C_p)^2 + (B_n - B_p)^2}$ 
end
k = GETK()
NED = RANK(ED)
PRESENT(NED[1 k])

```

Fig 2: pseudo code for Euclidean Distance

```

CLASSES = {'CLASS1', 'CLASS2', 'CLASS3', 'CLASS4', 'CLASS5'}
|C1 = 0, |C2 = 0, |C3 = 0, |C4 = 0, |C5 = 0
Dt ← Read(data);
ΣC = GETCOL(Dt) // get all cholesterol vector
ΣBp = GETBP(Dt) // get all bp vector
mcpoint = Max(C)
mbpoint = Max(Bp)
for each n in C
start:
-Cn < 40 & Bpn >= 70 & Bpn <= 90
C1 ↔ Dn
-Cn >= 40 & Cn < 200 & Bpn >= 70 & Bpn >= 90 & Bpn <= 120
C2 ↔ Dn
-Cn >= 200 & Cn < 239 & Bpn >= 70 & Bpn >= 100 & Bpn <= 140
C3 ↔ Dn
-Cn >= 40 & Cn < 200 & Bpn >= 70 & Bpn >= 90 & Bpn <= 120
C4 ↔ Dn
-Cn >= 240 & Cn < mcpoint & Bpn >= 70 & Bpn >= 140 & Bpn <= mbpoint
C5 ↔ Dn
end:

```

Fig 3: kNN algorithm

### 3.2 Decision Tree

The decision trees are broadly inquired about answer for characterization errands. For some sensible and pragmatic errands, the tree created by calculations is not fathomable to end client because of the extensive information size and many-sided quality of social information concerning traits. Numerous tree approaches have been worked out to deliver easier and more far reaching trees with great exactness in arrangement, rearrangements of tree has for the most part of auxiliary real concern in respect to precision furthermore no down to earth endeavour has been made to investigate the writing from the point of improvement. The structure that composes the ways to deal with tree improvement and close the methodologies inside this system. The fundamental point of this examination is to give the specialists a fine outline of tree disentanglement methodologies and understanding their combinational capabilities.[1][2][3][7] Our point is to foresee Risk element of the coronary illness utilizing the two Decision and Enhanced choice tree [9]. To assemble choice tree we considers c4.5 which takes various information as information. Part data is computed to locate the part of the hubs.  $S$  is the information set, and  $An$  is the arrangement of qualities, the condition beneath figures the data pick up for pair of trait  $(A_i, A_j)$  in set  $A$

$$\text{Split information}(s, A) = - \sum_{i=1}^k \frac{s_i}{s} \log_2 \frac{|s_i|}{s} \quad (2)$$

#### Algorithm:

```

INITIALIZATION
ΣA = 0, ΣW = 0, ΣC = 0
AA = 50 //AGE THRESHOLD
VALUE
AC = 200
AW = 90
L = [L1 L2 L3 L4]
L1 O "AGE"
DT <= READ(DATA)
W <= GETW(DT)
C <= GETC(DT)
FOR EACH N IN DT
A(N) < AA & A(N) > AA
START
W(N) < '70' & CN < '275'
At 50 & NW = DT(N)
W(N) < '70' & CN < '275'
Ac 50 & OW = DT(N)
W(N) < '90' & CN < '275'
Ac 50 & NW = DT(N)
W(N) < '90' & CN < '275'
Ac 50 & OW = DT(N)
END
FOR EACH L IN L
START
L1 <= "AGE"
L2 [1] <= At(N) < 50
L2 [2] <= Ac(N) > 50
L3 [1] <= At 50 & OW
L3 [2] <= At 50 & NW
L3 [3] <= Ac 50 & OW
L3 [4] <= Ac 50 & NW
END;
RAW ASSIGNMENT
L4 [1] <= COUNT [L3 [1]]? "Y"
L4 [2] <= COUNT [L3 [1]]? "N"
L4 [3] <= COUNT [L3 [2]]? "Y"
L4 [4] <= COUNT [L3 [2]]? "N"
L4 [5] <= COUNT [L3 [3]]? "Y"
L4 [6] <= COUNT [L3 [3]]? "N"
L4 [7] <= COUNT [L3 [4]]? "Y"
L4 [8] <= COUNT [L3 [4]]? "N"
TREE [L];

```

Fig 4: algorithm for decision tree

### 3.3 Improved Decision Tree

Improved Decision tree calculation discover nearby best answer for every Decision hub. To break out from nearby optima and to locate the worldwide arrangement we pick pair of properties at the same time, not one quality. Upgraded choice Tree technique, in picking qualities considers the data increase of picking a couple of characteristics simultaneously as opposed to picking one and only property. Thusly, to enhance the likelihood of result worldwide ideal arrangement, considering pair ideal characteristic is superior to anything single property. Upgraded Decision tree considers the pair of credits to develop leaf level hubs to diminish the quantity of levels.

$$\text{informationgain}(s, A_i, A_j) = \text{Entropy}(S) - \sum_{\substack{x \in \text{value}(A_i) \\ x \in \text{value}(A_j)}} \frac{S_{x,u}}{S} \text{Entropy}(S_{x,u}) \quad (3)$$

S is the information set, and An is the arrangement of properties, the condition beneath figures the data pick up for pair of trait (Ai, Aj) in set A. Data increase is utilized to locate the part of the pair of traits or a solitary characteristic.

#### Algorithm:

```

ΣG=0, ΣA=0, ΣW=0, ΣCP=0//CHEST PAIN
ΣCf=0 //CHOLESTEROL FOR FEMALE
ΣCm=0 //CHOLESTEROL FOR MALE
ΣS=0 //SMOKING
ΣA=0 //ALCOHOL
L=[L1 L2 L3 L4]
L16 " GENDER"
Dt<= READ(DATA)
G <= GET(Dt)
A<= GET(Dt)
for each n in Dt
start
L3[1]= GETM[S ,A ,CM ,0-50]
L3 [2]= GETF[W ,CP ,Cf ,0-50]
L3 [3]= GETM[S ,A ,CM ,50-100]
L3 [4]= GETF[W ,CP ,Cf ,50-100]
L4[1]= COUNT[L3[1], "Y"]
L4[2]= COUNT[L3[1], "N"]
L4[3]= COUNT[L3[2], "Y"]
L4[4]= COUNT[L3[2], "N"]
L4[5]= COUNT[L3[3], "Y"]
L4[6]= COUNT[L3[3], "N"]
L4[7]= COUNT[L3[4], "Y"]
L4[8]= COUNT[L3[4], "N"]
End TREE[L]

```

Fig5: pseudo code for enhanced decision tree

## IV. ANALYSIS

### 4.1 Precision

Accuracy alludes to the closeness of two or more estimations to each other. Accuracy (P) is characterized as the quantity of genuine positives (Tp) over the quantity of genuine positives included with the quantity of false positives (Fp). [3]

$$P = \frac{T_p}{T_p + F_p} \quad (4)$$

#### 4.2 Recall

The measure of genuinely significant results returned is called recall. Review (R) is characterized as the quantity of genuine positives (Tp) over the quantity of genuine positives in addition to the quantity of false negatives (Fn).[3]

$$R = \frac{T_p}{T_p + F_n} \quad (5)$$

#### 4.3 F-Measurement

F-measure is the measure of exactness test. It considers the accuracy p and the review r of the test to register the score: p is the quantity of right positive result isolated by the quantity of every single positive result, and r is the quantity of right positive results separated by the quantity of positive result that ought to returned. The F1 score is translated as a weighted normal of the accuracy and review, where the F1 score achieves the best esteem as 1 and most exceedingly terrible as 0.[3]

$$F = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} \quad (6)$$

Precision, Recall and F-measure are the precision estimations performed on the medicinal information to anticipate the precise calculation among Decision Tree, Enhanced choice tree and kNN

### V. OBSERVATIONS

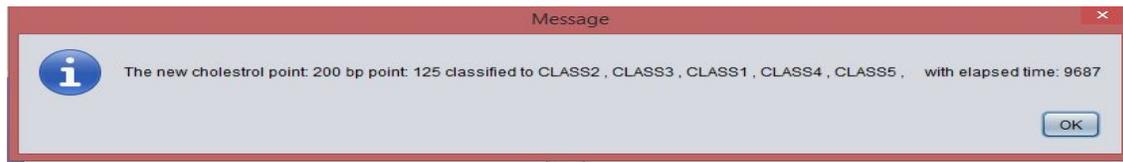


Fig 6: Class prediction using kNN

The kNN algorithm predicts the class of new attributes which are nearest to the input attributes with elapsed time.

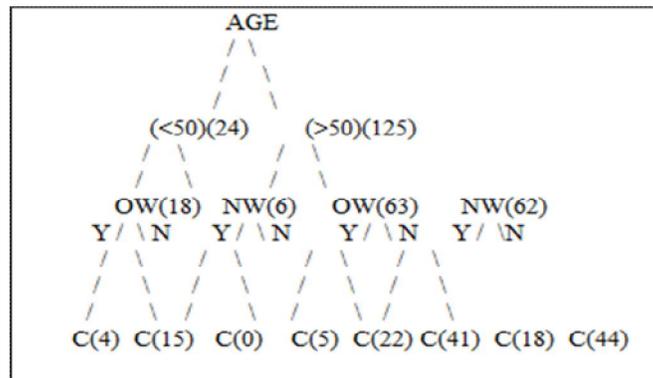


Fig 7: decision tree

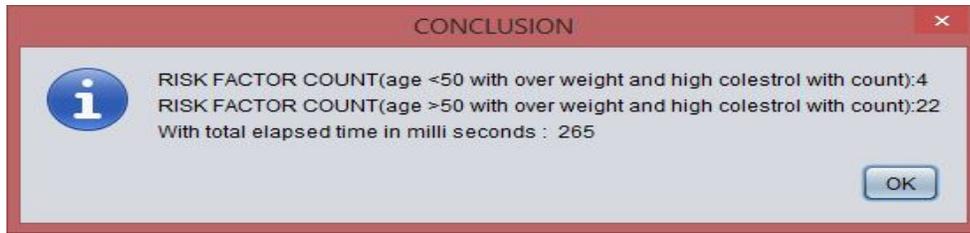


Fig 8: risk factor with elapsed time

Fig 7, 8 shows the results of the normal decision tree which predicts the Risk factor of heart disease with count of records and elapsed time.

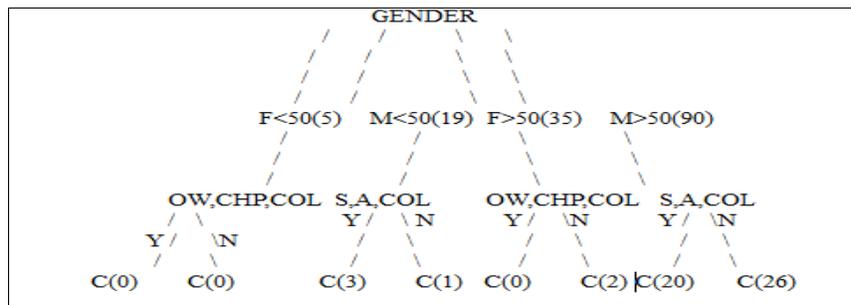


Fig 9: Improved decision Tree

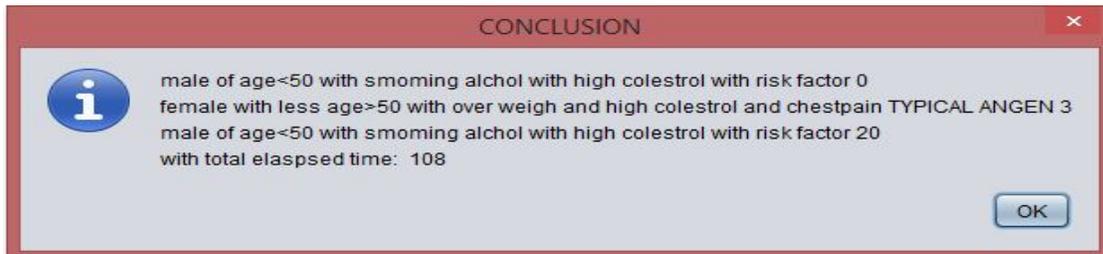


Fig 10: Risk factor with elapsed time

Fig 9, 10 demonstrates the consequence of Enhanced Decision tree with expectation of danger component and by considering more than two properties at leaf level and deliver the check of records which have the risk element.

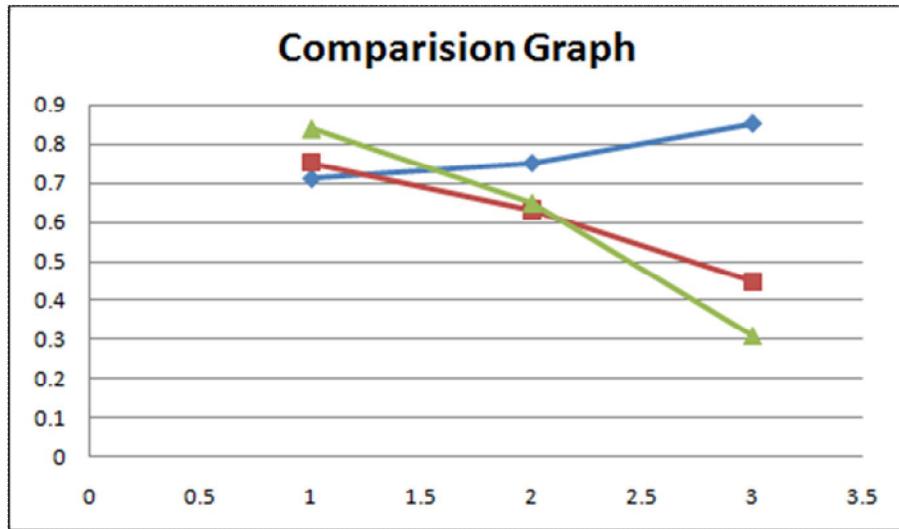


Fig 11: analysis graph

By the above diagram considering all estimations of exactness, review and time slipped by, of machine learning based kNN calculation, choice tree and Enhanced choice tree are appeared. Improved Decision tree delivers roughly exact results than other two methods on medicinal information.

## VI. CONCLUSION AND FUTURE WORK

This work can be stretched out to SAAS on W3C for web4.0 designs. Coming to general stubs and skeletons both sides the calculation part can be inundated taking into account augmentation of both information sets and methodologies relocation. The SOA design movement of this work is effortlessly supportive for further preparing of information like pre-handling and quick examination of different models and ascribed information. Tossing the information to cloud and on W3C of these sort of datasets constantly touchy however for security and analization reason stubs and skeletons keeps up security to prepare the information in mining and arrangement levels. For example disclosure this work can be stretched out from KNN level without further handling of choice trees. This is supportive in current work for different purposes.

## REFERENCES

- [1] Jiawei, H. (2006). Data Mining: Concepts and Techniques, Morgan Kaufmann publications.
- [2] Quinlan, J. R. (2014). C4. 5: programs for machine learning. Elsevier.
- [3] Karthikeyan, T., Thangaraju P. (2013). Analysis of Classification Algorithms Applied to Hepatitis Patients, International Journal of Computer Applications (0975 – 888), Vol. 62, No.15.
- [4] Suknovic, M., Delibasic B., et al. (2012). Reusable components in decision tree induction algorithms, Comput Stat, Vol. 27, 127-148.
- [5] Ruggieri, S. (2002). Efficient C4. 5 [classification algorithm]. Knowledge and Data Engineering, IEEE Transactions on, Vol. 14, No.2, 438-444.
- [6] Cios, K. J., Liu, N. (1992). A machine learning method for generation of a neural network architecture: A continuous ID3 algorithm. Neural Networks, IEEE Transactions on, Vol. 3, No.3, 280-291.
- [7] Gladwin, C. H. (1989). Ethnographic decision tree modeling Vol. 19.Sage.
- [8] Teach R. and Shortliffe E. (1981). An analysis of physician attitudes regarding computer-based clinical consultation systems. Computers and Biomedical Research, Vol. 14, 542-558.
- [9] Turkoglu I., Arslan A., Ilkay E. (2002). An expert system for diagnosis of the heart valve diseases. Expert Systems with Applications, Vol. 23, No.3, 229–236.
- [10] Witten I. H., Frank E. (2005). Data Mining, Practical Machine Learning Tools and Techniques, 2<sup>nd</sup> Elsevier.
- [11] Herron P. (2004). Machine Learning for Medical Decision Support: Evaluating Diagnostic
- [12] Performance of Machine Learning Classification Algorithms, INLS 110, Data Mining..