# Application of Rough Sets to Predict the Breast Cancer Risk Association with Routine Blood Analyses

**Amr H. Abdel Haliem, Mohammed A. Atiea, Mohammed E. Wahed, Mohammed S. Metwally**

*Abstract: For women around the globe, breast cancer has been a significant cause of mortality. Around the same time, early diagnosis and high cancer prediction precision are critical to improving the quality of care and the recovery rate of the patient. Expert systems and machine learning techniques are gaining prominence in this area as a result of efficient classification and high diagnostic ability. This paper introduces a model of hybrid prediction (RS QA) based on a rough set theoryand a quasi-optimal (AQ) rule induction algorithm. To find a minimal set of attributes that completely define the results, a rough set tool is used. The selected characteristics were collected, ensuring the high standard of the classification. Then to produce the decision rules, we use the quasi-optimal (AQ) rule induction algorithm. These hybrid prediction models allow expert systems to be built based on the conceptual rules of the IF/THEN sort. The suggested experiment is performed using the Coimbra Breast Cancer Dataset (BCCD) based on sets of measures that can be obtained in routine blood tests. Using classification precision, sensitivity, specificity, and receiver operating characteristics (ROC) curves, the efficiency of our suggested approach was assessed. Experimental results indicate the highest classification accuracy (91.7 percent), sensitivity (83.3 percent), and precision (94.3) obtained by the proposed (RS_QA) model.*

*Keywords—rough set, breast cancer, prediction system, AQ algorithm, rule induction.*

## I. INTRODUCTION

The body consists of several types of tissues that are made up of several cells, and in a healthy body, cells grow and divide into new cells and die in an orderly way. Cancer begins to appear while cells develop in inharmonious (*i.e.* out of control) [1,2]. Breast cancer is a disease that results from the uncontrolled growth of cells within the milk-producing glands in the breast Figure 1 highlights characteristics of most cancers at the site of the breast.Breast most cancers are the second leading component of cancer demise amongst American girls and the mortality of this disease decreased by using 40% between 1989 and 2016. The cause for this decline is early detection.
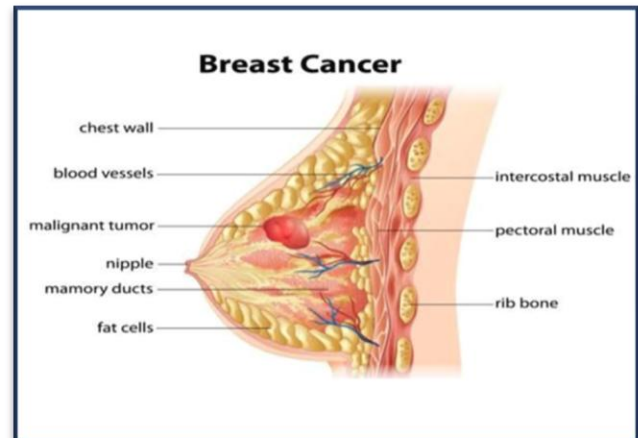


**Figure. 1. characteristics of cancer in the breast region.**

In science technology sickness analysis is a complex procedure, and a spread of exams is required for correct diagnosis. With the recent tendencies in artificial intelligence, machine learning techniques are applied to resource physicians to early diagnose a sickness reliably. Due particularly to initial detection of breast cancer, Patients should be supported with greater cure options, and loss of life threat may be minimized [3]. For breast most cancers wherein early analysis is essential, the usage of low fee and easy-to-apply methods like ordinary blood analysis is on hand for sufferers' lives earlier than mammography and magnetic resonance imaging (MRI).

The Rough Sets Theory (RS) is a mathematical method that extracts information from ambiguous, incomplete data. This theory states that we first have the facts or experience required to classify the objects into separate classes. When we have precisely the same details about two things, they are indiscernible (similar) i.e., we tell them we can't differentiate them by known knowledge [4]. RS theory can be used to find relationships of dependency between data, to calculate the significance of attributes, to reduce all redundant, and to look for the minimum attribute subsection to achieve satisfying classification. In terms of medicine, this is meant to define sub-sets of the most important features affecting medical care. Algorithm almost-optimum (AQ) is an algorithm of rule induction used to achieve meaningful rules of judgment to identify new cases and to expose deep medical information and provide new medical experience. These rules of decisions are more useful for assessing and interpreting the situation at hand by medical professionals. A rough set is a good methodology for assessing medical and health data. Hamouda et al.

**Amr H. AbdelHaliem \***, Faculty of science, computer science department, Suez, Egypt Email: amr_cs_2012@yahoo.com

**Mohammed A. Atiea**, Faculty of Computers and Informatics, Suez University, Suez, Egypt M_ail_atiea@hotmail.com

**Mohammed E.Wahed**, Faculty of Computers and Informatics, Suez University Ismailia, EgyptEmail: mewahed@yahoo.com

**Mohammed S.Metwally,** Faculty of Science, Department ofMathematics, Suez University, Suez, Egypt Email: met641958@yahoo.com

[5] suggested a breast cancer (BC) prediction and classification system classifying the 9-stage BC, which also sees the main attribute determining the BC stage. Inbarani [6] suggested a new approach to classification based on neighborhood rough set (NRS). The rough sets approach used by Wang et al. [7] to predict brain glioma malignancy with dataset contain irrelevant features and missing values. Fakih et al. [8] extract diagnostic rules using rough sets from the medical databases. Michalowski et al. [9] introduced a rough approach to select the most important clinical characteristics and to create decision rules from medical databases with missing values. Hassanien [10] To analyze data on breast cancer utilizing rough-set theory. Chul-Heui Lee et al [11] suggested using the rough set theory for a new classification system.

In this paper, we are presenting a Hybrid Prediction System (RS_QA) based on a rough set and quasi-optimal (AQ) rule induction algorithm to predict breast cancer based on routine blood analysis features. A rough set algorithm is used to select a more powerful subset of features. The selected subsets are then used in the decision rule generation process to create descriptive rules for the classification task Using the quasi-optimal (AQ) rule induction algorithm.

The paper is planned as follows. in Section 2 The features of breast cancer data are considered and essential knowledge about rough sets approach and rule induction are reviewed, which are related to the work, and the AQ (rule induction algorithms). Sections 3 present the proposed Prediction System (RS_QA). The results of these experiments are demonstrated in Section 4. the conclusion is exposed in Section 5.

## II.  MATERIALS AND METHODS

### A.  The characteristics of breast cancer data

In this paper, the proposed approach was evaluated for disease diagnosis on Breast Cancer Coimbra Dataset collected via routine blood analysis, publicly available on the University of California–Irvine Machine Learning Repository [12]. The dataset includes 116 instances in which 64 of them are breast cancer patients (i.e., positive class) and 52 of them are healthy controls (i.e., negative class). There are nine quantitative attributes, and one binary class attribute indicating the patients or healthy controls. The overview of the Breast Cancer Coimbra Dataset containing attribute is shown in Table 1.

**TABLE I.        : SUMMARY OF ATTRIBUTES FOR BREAST CANCER COIMBRA DATASET**

| S No. | Attributes | UNIT |
|---|---|---|
| 1 | Age | years |
| 2 | BMI | kilograms per meter square (Kg/m$^2$) |
| 3 | Glucose | milligrams per decilitre (mg/dL) |
| 4 | Insulin | microunit per milliliter (µU/mL) |
| 5 | HOMA | ----- |
| 6 | Leptin | nanogram  per milliliter (ng/mL) |
| 7 | Adiponectin | microgram per milliliter(µg/mL) |
| 8 | Resistin | nanogram per milliliter (ng/mL) |
| 9 | MCP-1 | picogram per decilitre( pg/dL) |

### B.  Rough Set Theory

As an effective technique for overcoming uncertainty and confusion in data analysis, Pawlak [13] introduced the

principle of Rough Sets. For example, if objects are patients with a certain illness the symptoms are information about the patient. The Objects thatlabeled by the same attribute value (similar) are indiscernible by the knowledge accessible. The relation of indiscernibility that this prompt is the basis of a mathematical definition. A main recommendation within a roughly defined learning system is to recognize redundancies and dependencies between the features of a problem that should be classified. It depends on the concept of the upper and lower approximation of the set, the approximation of the space, and the models of the set. The primary benefit of RS theory is that there is no need for external knowledge on data such as statistical probability or membership grade, as in fuzzy set theory.

Consider that $U$ is an object set and $A$ is an attribute set. Suppose that $U$ and $A$ are finite sets.Then the pair $(U, A)$ is called an information system $(IS)$, if each attribute $a \in A$ determines an information function $f_a$: $U \rightarrow Va$, where $Va= \{a(x): x \in U\}$ is the set of information function values of the attribute $a$. If $(U, A)$ is an $IS$ and $A = C \cup \{d\}$ where $C$ is a conditional attribute set and $d$ is a decision attribute set. Then $(U, C \cup \{d\})$ is referred to as a decision system.

For any information system, we can express a relation between objects with their attribute values. For a given attribute $a$; the objects $x$; $y$are a-indiscernible if they have the same value on $a$, i.e. $a(x) = a(y)$.In these terms, we call two objects indiscernible if one cannot distinguish between them using only the knowledge available in the decision table. This definition can be presented to any subset B $\subset$ A by:

$$IND(B) = \{(x, y) \in U \times U | \forall a \in B, f_a(x) = f_a(y)\} \qquad (1)$$

IND(B) denotes the relation determined on the subset B $\subset$ A. IND(B) is an equivalence relation. Objects x; y are indiscernible by attributes from B. We denote by $[x]_B$ The equivalence classes defined by the object x $\in$ U.

For any $X \subset U$, one can define the lower approximation and the upper approximation of X by:

$$\underline{B}(X) = \{x \in U \mid [x]_B \subseteq X\} \text{ lower approximation}$$

$$\overline{B}(X) = \{x \in U \mid [x]_B \cap X \neq \phi\} \text{ upper approximation}$$

The upperapproximation consists of data points that are possibly (plausibly) similar to those inX according to the knowledge taken from the decision table. The lower approximation of X is a set of data points that are with certainty similar only to the elements ofX [14,15].

the universe $U$ can be divided into three regions (positive, boundary, and negative regions) as follows [13,14]:

$POS_B(A) = \underline{B}(x)$ denotes the positive region of $X$,

$NEG_B(X) = U - \overline{B}(X)$ denotes the negative region of $X$ and $Bd_B(X) = \overline{B}(X) - \underline{B}(X)$ denotes the boundary region are shown in Figure. 2
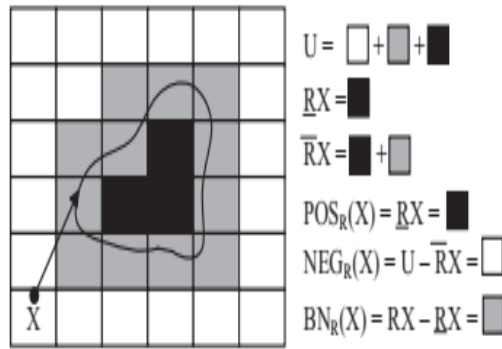
**Figure 2: The approximations and regions of set X using Rough Set theory.**

The identification of a reliance on the value range of attributes in a data analysis Rough setting is an important problem and can be numerically calculated as $\alpha_B(X)$

$$\alpha_B(X) = \left| \frac{\underline{B}(X)}{\overline{B}(X)} \right|$$

Reduce and are central principles in rough set theory. A reduct is a subset of attributes such that it is enough to consider only the features that belong to this subset and still have the same amount of information. Moreover, the reduct has the property of minimality i.e. it cannot be reduced any more without loss in quality of information. there can be several reducts for a decision table (information system) and *RED(A)* can denote that as a family of reducts for a decision table. the core can represent the least possible number of attributes *(*intersect*).*

$$CORE(A) = \cap RED(A)$$

And represent the part of information that cannot be removed from the system without loss in knowledge that can be derived from it.

### C. Rule Induction

#### 1) Rule-based algorithms

The rule-based induction strategy for data mining to find the underlying concepts in the data set as the regularity of the rules.

The rules are usually formulas:

if (attribute − 1, value − 1) and (attribute − 2, value − 2) and ⋯and (attribute − n, value − n) then (decision, value).

*Algorithm quasi-optimal (AQ) is a* rule induction algorithm, developed by R. S. Michalski [16] [17]. defined as the supervised learning algorithm. AQ establishes rules through an iterative mechanism to identify inferences of positive examples about negative examples (i.e., rules from examples and counterexamples are generated).

Algorithm 1 showsthe AQ algorithm, where P represents positive examples (cancer patient in our case), and N represents negative examples (normal in our case).the algorithm begins by randomly selecting seed from positive event list P*, and then creates a STAR rule for that example, which is an iterative process aimed at generating a range of possible descriptions (rules) of the seed that satisfy certain limitations, for example, do not cover negative examples, do not contradict prior knowledge. By first classifying, then generalizing each new example to the best previously generated rule set, the AQ algorithm learns incrementally. Before a new example is generalized by AQ, it checks out if

there are any rule sets in the affected feature space area that conflict with the new rule set proposed and with the opposite class. If so, the generalization is aborted and exactly storage of the record is done.

| Algorithm 1: AQ |
|---|
| Input : |P| > 0 & |N| > 0 |
| Output: C |
| 1. P* = P; C = 0 //P* is a list of positive events to be covered |
| 2. **While** |P*| > 1 **do** |
| 3. Select random seed *p* from P* |
| 4. c = STAR (*p*, N, maxstar) //find a rule that generalize the seed |
| 5. P* = P* - [ P* ∩ r]   // Remove from P examples covered by c |
| 6. C = C + c //increment the set of rules by new one |
| 7. **end while** |
| Return C |

#### 2) Classification accuracy

Using Eq 2, the accuracy of the values obtained is determined. Accuracy is classified as the rate of correctly classifying data that was not previously found by the algorithm.

$$Accuracy(\%) = \frac{No.of\ test\ examples\ covered\ by\ the\ rule\ set}{Total\ no.of\ test\ examples} \quad (2)$$

The classification accuracy determines as the count of tests correctly classified by the rulesets, divided by the size of the test set chosen from the sample set. The mean average of the values calculated according to the 10-fold cross-validation method. another common metric for the validation of classification algorithms in this paper. are as follows:

$$Specficity = \frac{TN}{FP + TN}$$

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

$$AUC = \frac{1}{2}(Sensitivity + specificity)$$

where, TP = true positive, TN = true negative, FN = false-negative and FP = false positive.

## III. METHODOLOGY

*Proposed methodological flowchart*

First step preprocessing data that includes discretization. Then we will use concepts of rough sets theory as a Feature selection technique to reduce the superfluous attributes. The second step using the AQ algorithm to extract a minimum set of rules to use for the classification process. the flowchart in Fig 3 shows that.

## IV. EXPERIMENTS RESULTS AND DISCUSSION

In our studies, data preprocessing works a major role in the generation of rules. First, we calculated the indiscernibility relation for the dataset. Then determined the approximations sets (upper, lower) to evaluate the rough sets. Finally, the positive region and the degree of dependence have been determined.

# Application of Rough Sets to Predict the Breast Cancer Risk Association with Routine Blood Analyses

The dataset has evidence of the existence rough set because the calculated dependence is 0.976. After the data preprocessing stage, the decision rules are generated by the AQ algorithm. There is a series of 105 rules extracted using this algorithm to predict whether or not breast cancer has been observed.

## A. Rule-based – classification

The accuracy of classification of the set of rules derived from the data set shall be checked by 10-fold cross-validation. All cases are randomly divided into ten folds of equal size at random. For rule induction (training), using each complete data subset, while the remaining subset is utilized for testing. We execute 10 times on each test and the results are averaged. The ROC curve plots the sensitivity (percentage of correctly classifying positive subjects) versus the specificity or (percentage of incorrectly classifying negative class) for every classification subset. the AUC (area under the ROC curve) is computed to describes the global precision.Its detailed information on Evaluation Measures is shown in Table 2. And Figures 2 ROC curve and results of area under the curve (AUC) of our proposed approach. From Figure4, we can see that the results of our proposed approach can achieve promising results for breast cancer diagnosis. From Figure4, we can see that the AUC of BCCD datasets are 0.91.

### TABLE II. THE DETAILED INFORMATION ABOUT EVALUATION MEASURES

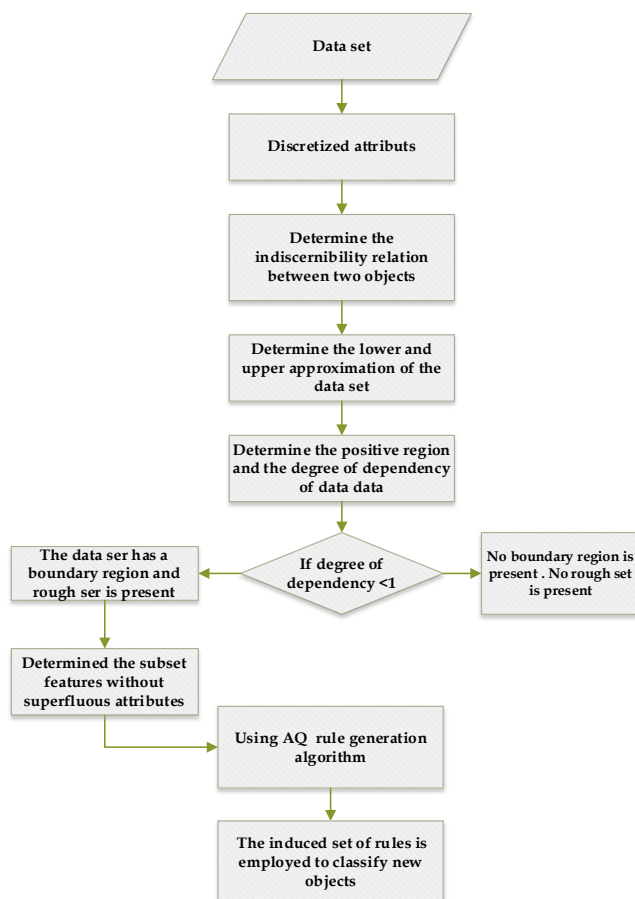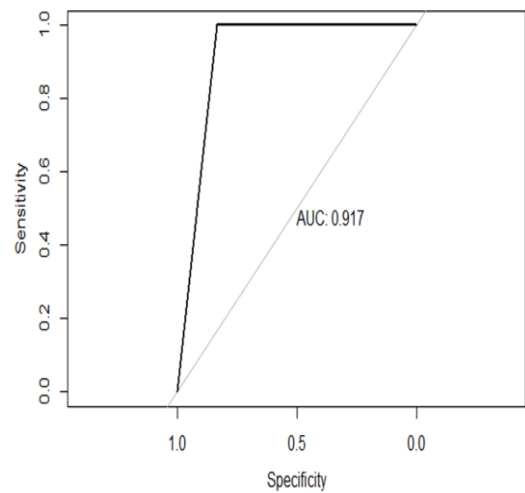| Proposed classification model | accuracy | sensitivity | specificity |
|---|---|---|---|
| | 0.91 | 0.833 | 0.9412 |



**Figure 3: proposed framework**



**Figure 4: The accuracy of the proposed framework (RST-AQ)**

## B. Decision rules generated from breast cancer data

There is a set of 105 certain rules extracted by using this algorithm. Only 10 rules for Breast cancer (cancer) cases are shown and explained in Table 3. For example, Rule 1, IF BMI is (25.1,30.4] and HOMA is (1.85, Inf] and Resistin is (14.2, Inf], covers 12 of 64 breast cancer cases. Rule 2, IF Glucose is (96.3, Inf] and Resistin is (14.2, Inf], covers 11 of 64 breast cancer casesAlso, IF BMI is (25.1,30.4] and Adiponectin is (6.68,10.6) and Resistin is [-Inf,8.47] always leads to breast cancer. This rule covers 7 of 64 breast cancer cases. From these rules, the following conclusions can be drawn: If (BMI is (25.1,30.4]) AND (Resistin is (14.2, Inf]) AND (Glucose is (96.3, Inf]) AND (HOMA is (1.85, Inf]) AND (MCP.1 is [-Inf,334]) Then (certain cases will be cancer)

## C. Performance comparison with other algorithms

The proposed approach (RS -AQ) was compared with existing works using the same dataset (BCCD) with different methods. Table 4 presents a comparison list for accuracy metrics.

### TABLE III. ONLY 10 RULES GENERATED FROM BREAST CANCER DATA

| rule | If condition | Then condition | Number of elements supporting the condition |
|---|---|---|---|
| R1 | IF BMI is (25.1,30.4] and HOMA is (1.85, Inf] and Resistin is (14.2, Inf] | C | 11 |
| R2 | IF Glucose is (96.3, Inf] and Resistin is (14.2, Inf] | C | 11 |
| R3 | IF BMI is (25.1,30.4] and Adiponectin is (6.68,10.6) and Resistin is [-Inf,8.47] | C | 7 |
| R4 | IF Resistin is (14.2, Inf] and Leptin is (14.8,29.8] and MCP.1 is (619, Inf] | C | 6 |
| R5 | IF Glucose is (96.3, Inf] and Age [-Inf,47.7] | C | 6 |
| R6 | IF Insulin is (4.52,7.73] and Glucose is [-Inf,87] | C | 12 |

| R7 | IF BMI is [-Inf,25.1] and MCP.1 is [-Inf,334] and Glucose is (87,96.3] | C | 5 |
|---|---|---|---|
| R8 | IF MCP.1 is [-Inf,334] and Resistin is (14.2, Inf) | C | 5 |
| R9 | IF Glucose is (87,96.3] and BMI is [-Inf,25.1] and MCP.1 is [-Inf,334] | C | 5 |
| R10 | IF Insulin IS (7.73, Inf] and Glucose IS (96.3, Inf] and Leptin is (29.8, Inf] | C | 11 |

**TABLE IV. PERFORMANCE METRICS OBTAINED WITH OUR PROPOSED FRAMEWORK AND OTHER METHOD OBTAINED FROM THE LITERATURE FOR THE BREAST CANCER COIMBRA DATASET.**

| Algorithm | method | Average accuracy | Average Sensitivity | Average Specificity |
|---|---|---|---|---|
| M.Patrício [18] | SVM | [0.86, 0.90] | 0.84 | [0.81,0.87] |
| | logistic regression | [0.79, 0.83] | 0.79 | [0.81, 0.87] |
| | random forests | [0.82, 0.87] | 0.85 | [0.77, 0.83] |
| Çelik,Yunus,et al [19] | Artificial Neural Network (ANN) | 79.4 | -- | -- |
| | Standard Extreme Learning Machine (ELM) | 80.0 | -- | -- |
| | Support Vector Machine (SVM) | 77.5 | -- | -- |
| | K-Nearest Neighbor(k-NN) | 73.5 | -- | -- |
| Silva Araújo et al[20] | hybrid system combining fuzzy systems and neural networks | 0.81 | | |
| Akben[21] | decision tree | 0.90 | | |
| Singh[22] | medium Gaussian SVM classifier | 0.82 | | |
| Li and Chen[23] | (ANN), decision tree, SVM, random forests, and logistic regression | 0.74 | | |
| Proposed framework | Rough set + AQ | **0.917** | **0.833** | **0.9412** |

## V. CONCLUSIONS

Experimental results showed that the proposed algorithm was more efficient and generated decision rules with better classification performance. Features such as Resistin, Glucose, BMI, MCP.1 are crucial to the degree prediction of breast cancer proposed framework has effectivelycombined the basic notions of rough sets theory with AQ rule induction algorithm to yield a new approach for inductive learning under uncertainty. Through this integration, mixed the potencies of rough-sets theory and AQ rule induction algorithm.

The opportunity to form simple and clear rules offers the following benefits:

- It is easy to understand.
- It is easy to interpret and analyses.
- It is easy to check and verify.

For medical applications, rough set feature selection and rule induction approaches are powerful to interpret medical data even though complexity and missing values occur.

## REFERENCES

1. National breast cancer foundation, (2017), Available online: http://www.nationalbreastcancer.org/about-breast-cancer. (accessed on 10 January 2020).
2. Cancer.net editorial board, (2017), Available online: http://www.cancer.net/ cancer-types/breast-cancer. (accessed on 10 January 2020).
3. Jemal A, Murray T, Ward E, Samuels A, Tiwari RC, Ghafoor A, Feuer EJ, Thun MJ. Cancer statistics, 2005. CA: a cancer journal for clinicians. 2005 Jan 1;55(1):10-30.
4. Tseng, T. L.(Bill), Quantitative approaches for Information modeling, Ph.D. Dissertation, University of Iowa.
5. S.K.M. Hamouda, M.E. Wahed, R.H. AboAlez, and K. Riad. Robust breast cancer prediction system based on rough set theory at the national cancer institute of Egypt. Computer Methods and Programs in Biomedicine, 153,2018:259–268.
6. Inbarani HH. A novel neighborhood rough set-based classification approach for medical diagnosis. Proc Comput Sci 2015;47:351–9.
7. Wang, X., Yang, J., Jensen, R., & Liu, X. Rough set feature selection and rule induction for prediction of malignancy degree in brain glioma. Computer Methods and Programs in Biomedicine, 83(2), (2006), 147–156.
8. Fakih, S. J., & Das, T. K. LEAD: A methodology for learning efficient approaches to medical diagnosis. IEEE Transactions on Information Technology in Biomedicine, 10(2), (2006), ppt.220–228.
9. Wilk, Sz., Slowinski, R., Michalowski, W., & Greco, S. Supporting triage of children with abdominal pain in the emergency room. European Journal of Operational Research, 160(3) (2005), pp.696–709.
10. A.E. Hassanien, Rough set approach for attribute reduction and rule generation: a case of patients with suspected breast cancer, J. Am. Soc. Inform. Sci. Technol. 55 (11) (2004), pp. 954–962.
11. L.H. Chul, S.H. Seon, C.C. Sang, Rule discovery using hierarchical classification structure with rough sets, in: FSA World Congress and 20th NAFIPS International Conference, 1, 2001, pp. 447–452.
12. UCI. "Machine Learning Repository," https://archive.ics.uci.edu/ml/index.php. . (accessed on 10 January 2020).
13. Z. Pawlak," Rough Sets", International Journal of Computer and Information Sciences, vol. 11, 1982, pp. 341-356.
14. Z. Pawlak, "Rough Sets Theoretical Aspects of Reasoning about Data", Kluwer Academic Publishers, Dordrecht, 1991.
15. J. Komorowski," Learning Rule-Based Models The Rough Set Approach", Comprehensive Biomedical Physics, vol. 6, 2014,pp. 19-39.
16. Michalski R.S., Mozetic I., Hong J., Lavrac N. The AQ15 inductive learning system: An overview and experiments, Report 1260, Department of Computer Science, University of Illinois at Urbana-Champaign, 1986A.
17. Michalski R.S., Mozetic I., Hong J., Lavrac N. The multi-purpose incremental learning system AQ 15 and it is testing application to three medical domains. Proc. of the 5th Nat. Conf. on AI, 1986B, 1041-1045.
18. Patrício, Miguel , Pereira, José , Crisóstomo Silva, Joana , Matafome, Paulo , Gomes, Manuel , Seiça, Raquel , Caramelo, Francisco. (2018). Using Resistin, glucose, age and BMI to predict the presence of breast cancer. BMC Cancer. 18. 10.1186/s12885-017-3877-1.
19. Çelik, Yunus & Sabanci, Kadir & Durdu, Akif & Aslan, MuhammetBreast Cancer Diagnosis by Different Machine Learning Methods Using Blood Analysis Data. International Journal of Intelligent Systems and Applications in Engineering. 6, (2018), pp. 289-29
20. Silva Araújo VJ, Guimarães AJ, de Campos Souza PV, Silva Rezende T, Souza Araújo V, Using resistin, glucose, ageand BMI and pruning fuzzy neural network for the constructionof expert systems in the prediction of breast cancer. Mach Learn Knowl Extr 1(1): (2019),466–482

21. Akben SB,Determination of the bloodhormone and obesity value ranges that indicate the breast cancer, using data mining based expert system. IRBM 40(6): (2019) 355–360.

22. Singh BK,Determining relevant biomarkers for prediction of breast cancer using anthropometric and clinical features: a comparative investigation in machine learning paradigm. Biocybern Biomed Eng 39(2): (2019) 393–409.

23. Li Y, Chen Z, Performance evaluation of machine learning methods for breast cancer prediction. Appl Comput Math 7(4): (2018) 212-216.

## AUTHORS PROFILE

**Amr H. Abdel Halim** was born in 1985. He received B.Sc. in Mathematics and Computer Science in 2006 from Zagazig University, Zagazig, and M.sc in Computer Science in the same university in 2010. His area of interest includes developing data mining and soft computing teaching in medical applications**.**

**Mohammed E. Wahed** Doctor of Philosophy in Science (Mathematics and Operations Research and Computing) on 1995 from Zagazig University Professor of Computer Science, Department of Computer Science, Faculty of Computing and Information, Ismailia, Suez Canal University.His area of interest includes computer science, mathematics, operations research and decision support system.