



# Analysis of Breast Cancer Detection Using Different Machine Learning Techniques

Siham A. Mohammed<sup>1</sup>, Sadeq Darrab<sup>3</sup>(✉), Salah A. Noaman<sup>2</sup>, and Gunter Saake<sup>3</sup>

<sup>1</sup> Taiz University, Taiz, Yemen

siham.alkadasi@gmail.com

<sup>2</sup> Aden University, Aden, Yemen

s-salah-17@hotmail.com

<sup>3</sup> University of Magdeburg, Magdeburg, Germany

{sadeq.darrab, gunter.saake}@ovgu.de

**Abstract.** Data mining algorithms play an important role in the prediction of early-stage breast cancer. In this paper, we propose an approach that improves the accuracy and enhances the performance of three different classifiers: Decision Tree (J48), Naïve Bayes (NB), and Sequential Minimal Optimization (SMO). We also validate and compare the classifiers on two benchmark datasets: Wisconsin Breast Cancer (WBC) and Breast Cancer dataset. Data with imbalanced classes are a big problem in the classification phase since the probability of instances belonging to the majority class is significantly high, the algorithms are much more likely to classify new observations to the majority class. We address such problem in this work. We use the data level approach which consists of resampling the data in order to mitigate the effect caused by class imbalance. For evaluation, 10 fold cross-validation is performed. The efficiency of each classifier is assessed in terms of true positive, false positive, Roc curve, standard deviation (Std), and accuracy (AC). Experiments show that using a resample filter enhances the classifier's performance where SMO outperforms others in the WBC dataset and J48 is superior to others in the Breast Cancer dataset.

**Keywords:** Breast cancer · Classification · Data mining

## 1 Introduction

Breast cancer is the second leading cause of death among women worldwide [1]. In 2019, 268,600 new cases of invasive breast cancer were expected to be diagnosed in women in the U.S., along with 62,930 new cases of non-invasive breast cancer [2]. Early detection is the best way to increase the chance of treatment and survivability. Data mining has become a popular tool for knowledge discovery which shows good results in marketing, social science, finance and medicine [19, 20]. Recently, multiple classifiers algorithms are applied on medical datasets to perform predictive analysis about patients and their medical diagnosis [6, 9, 10, 21]. For example, using machine learning techniques to assess tumor behavior for breast cancer patients. One problem is that there is a class

imbalance in the training data, since the probability of not having this disease is higher than the one of having it. This paper introduces a comparison between three different classifiers: J48, NB, and SMO with respect to accuracy in detection of breast cancer. Our aim is to prepare the dataset by proposing a suitable method that can manage the imbalanced dataset and the missing values, to enhance the classifier's performance. All tasks were conducted using Weka 3.8.3.

The remainder of this paper is organized as follows. Section 2 presents literature review. Section 3 introduces the datasets. Section 4 describes the research methodology including pre-processing experiments, classification and performance evaluation criteria. The experimental results are presented in Sect. 5. Finally, Sect. 6 shows the conclusion and future work.

## 2 Literature Review

In recent years, several studies have applied data mining algorithms on different medical datasets to classify Breast Cancer. These algorithms show good classification results and encourage many researchers to apply these kind of algorithms to solve challenging tasks. In [21], a convolutional neural network (CNN) was used to predict and classify the invasive ductal carcinoma in breast histology images with an accuracy of almost 88%. Moreover, data mining is used widely in medical fields to predict and classify abnormal events to create a better understanding of any incurable diseases such as cancer. The outcomes of using data mining in classification are promising for breast cancer detection. Therefore, data mining approach is used in this work. A list of some literature studies related to this method is presented in Table 1.

## 3 Datasets

The datasets that are used in this paper are available at the UCI Machine Learning Repository [13].

### 3.1 WBC Dataset

The WBC dataset contains 699 instances and 11 attributes in which 458 were benign and 241 were malignant cases [14]. In the WBC, the value of the attribute (Bare Nuclei) status was missing for 16 records. Hence data preprocessing is essential and important for this dataset, requiring us to manage the imbalanced data and the missing values.

### 3.2 Breast Cancer Dataset

The feature form this dataset are computed from a digitized image of a fine needle aspirate (FNA) of a breast tumor. The target feature records the prognosis (i.e., malignant or benign). The dataset contains 286 instances and 10 attributes in which 201 were no-recurrence-events and 85 were recurrence events. In the Breast Cancer dataset, the value of the attribute (node-caps) status was missing in 8 records.

**Table 1.** Breast cancer detection research using different machine learning algorithms.

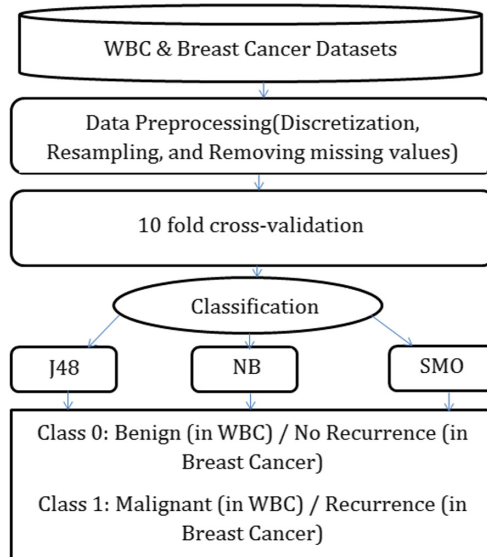
Paper title	Datasets	Algorithms	Results
Integration of data mining classification techniques and ensemble learning for predicting the type of breast cancer recurrence [3], 2019	Breast Cancer	NB, SVM, GRNN and J48	GRNN & J48 accuracy: 91% NB & SVM: 89%
A study on prediction of breast cancer recurrence using data mining techniques [4], 2017	WPBC	Classification: KNN, SVM, NB and C5.0, Clustering: K-means, EM, PAM and Fuzzy c-means	Classification accuracy is better than clustering, SVM & C5.0: 81%
Predicting breast cancer recurrence using effective classification and feature selection technique [5], 2016	WPBM	NB, C4.5, SVM	NB: 67.17%, C4.5: 73.73%, SVM: 75.75%
Using machine learning algorithms for breast cancer risk prediction and diagnosis [6], 2016	WBC	SVM, C4.5, NB, KNN	SVM outperform others: 97.13%
Study and analysis of breast cancer cell detection using Naïve Bayes, SVM and ensemble algorithms [7], 2016	WDBC	NB, SVM, Ensemble	SVM: 98.5%, NB & Ensemble: 97.3%
Analysis of Wisconsin breast cancer dataset and machine learning for breast cancer detection [8], 2015	WDBC	NB, J48	NB: 97.51%, J48: 96.5%
Comparative study on different classification techniques for breast cancer dataset [9], 2014	Breast Cancer	J48, MLP, rough set	J48: 79.97%, MLP: 75.35%, rough set: 71.36%
A novel approach for breast cancer detection using data mining techniques [10], 2014	WBC	SMO, IBK, BF Tree	SMO: 96.19%, IBK: 95.90%, BF Tree: 95.46%
Experiment comparison of classification for breast cancer diagnosis [11], 2012	WBC WDBC WPBC	J48, SMO, MLP, NB, IBK	In WBC: MLP & J48: 97.2818%. In WDBC: SMO: 97.7% or fusion on SMO & MLP: 97.7% In WPBC: fusion of MLP, J48, SMO and IBK: 77%
Analysis of feature selection with classification: breast cancer datasets [12], 2011	WBC WDBC Breast Cancer	Decision Tree with and without feature selection	Feature selection enhances the results WBC: 96.99% WDBC: 94.77% Breast Cancer: 71.32%

## 4 Research Methodology

The two datasets used in this work are vulnerable to missing and imbalanced data therefore, before performing the experiments, a large fraction of this work will be for pre-processing the data in order to enhance the classifier's performance. Preprocessing will focus on managing the missing values and the imbalanced data. To manage the missing attributes, all the instances with missing values are removed. The imbalance data problem needs to adjust either the classifier or the training set balance. To do so, the resample filter is used to rebalance the data artificially. Then, 10 fold cross validation is applied and finally a comparison between these three classifiers is implemented.

### 4.1 Preprocessing Phase

First, the data were discretized using discretize filter, then missing values were removed from the dataset. Second, instances were resampled using the resample filter in order to maintain the class distribution in the subsample and to bias the class distribution toward a uniform distribution. Section 5 will show that this idea is improving the classifier's performance. Third, 10 fold cross validation was applied then experiments were carried out over three classifiers Naïve Bayes, SMO and J48, as illustrated in Fig. 1.



**Fig. 1.** Proposed breast cancer detection model using Breast Cancer and WBC datasets.

In Fig. 1, the data preprocessing technique has been applied including three steps: discretization, instances resampling and removing the missing values. After that, 10 fold cross validation has been applied. Then, three classifiers have been evaluated over the prepared datasets.

## 4.2 Training and Classification

In order to minimize the bias associated with the random sampling of the training data, we use 10 fold cross validation after the pre-processing phase. In k-fold cross-validation, the original dataset is randomly partitioned into k equal size subsets. The classification model is trained and tested k times. Each time, a single subset is retained as the validation data for testing the model, and the remaining  $k-1$  subsets are used as training data. Three classification techniques were selected: a Naïve Bayes (NB), a Decision Tree built on the J48 algorithm, and a Sequential Minimal Optimization (SMO). The NB classifier is a probabilistic classifier based on the Bayes rule. It works by estimating the portability of each class value that a given instance belongs to that class [15]. The J48 algorithm [16] uses the concept of information entropy and works by splitting each data attributers into smaller datasets in order to examine entropy differences. It is an improved and enhanced version of C4.5 [17]. The SMO model implements John Platt's sequential minimal optimization algorithm for training a support vector classifiers. This implementation globally replaces all missing values and transforms nominal attributes into binary ones. It also normalizes all attributes by default [18].

## 4.3 Performance Evaluation Criteria

In this study, we use five performance measures to evaluate all the classifiers: true positive, false positive, ROC curve, standard deviation (Std) and accuracy (AC).

$$AC = (TP + TN)/(TP + TN + FP + FN). \quad (1)$$

Where TP, TN, FP and FN denote true positive, true negative, false positive and false negative, respectively.

# 5 Experimental Results

First, the three classifications algorithms were tested on the WBC and the Breast Cancer datasets without applying the preprocessing techniques. Among them, the best result was recorded for J48: 75.52% in the Breast Cancer dataset and for SMO: 96.99% in the WBC dataset. Next, after applying preprocessing techniques accuracy increases to 98.20% with J48 in the Breast Cancer dataset and 99.56% with SMO in the WBC dataset.

## 5.1 Experiment Using the Breast Cancer Dataset

First, the three classifiers are tested over original data (without any preprocessing). The results show that J48 is the best one with 75.52% accuracy where the accuracy of NB and SMO are 71.67% and 69.58%, respectively. Next, we apply discretization filter and remove the records with missing values, results improved with NB and SMO as follows: NB: 75.53% and SMO: 72.66% where J48: 74.82%. After that, resample filter was applied for 7 times. The Performance of the classifiers are improved and enhanced as shown in Table 2.

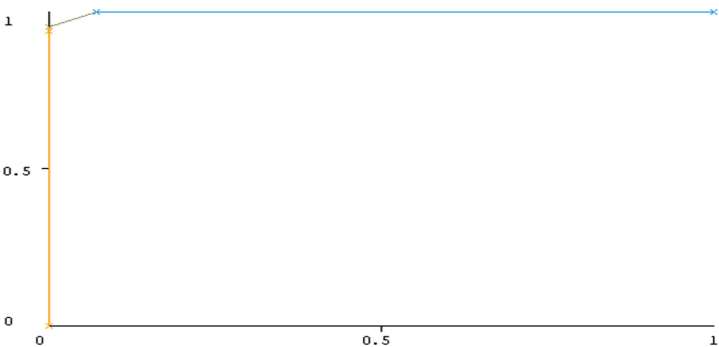
**Table 2.** Performance of the classifiers in the Breast Cancer Dataset.

Experiments steps	Classifier accuracy		
	J48	NB	SMO
Original without preprocessing	75.52%	71.67%	69.58%
After removing missing values & discretization	74.82%	75.53%	72.66%
After applying resample filter (first time)	79.49%	77.33%	80.93%
Applying resample filter (second time)	81.65%	78.05%	80.57%
Applying resample filter (third time)	87.41%	78.41%	82.73%
Applying resample filter (fourth time)	92.08%	77.69%	88.84%
Applying resample filter (fifth time)	95.68%	79.13%	91.72%
Applying resample filter (sixth time)	97.48%	79.85%	95.68%
Applying resample filter (seventh time)	98.20%	76.61%	95.32%

As illustrated in Table 2, we can obviously notice that the more resample filter we apply, the improved accuracy we obtain. That is because the data is imbalanced and the filter maintains the class distribution. For the Breast cancer dataset, J48 outperforms others with 98.20%. Accuracy measures for J48 classifier is shown in Table 3 and Roc curve of J48 is shown in Fig. 2.

**Table 3.** Accuracy measures for J48 in the Breast Cancer Dataset.

TP	FP	Precision	Recall	Roc curve	Std	Class
1.000	0.049	0.980	0.996	1.000	0.5678	No-recurrence-events
0.951	0.000	1.000	0.996	0.951		Recurrence-events



**Fig. 2.** J48 ROC curve in Breast Cancer Dataset.

To measure the performance of the proposed model, we compare the obtained results with the study proposed in [9]. The same dataset and three classifiers including J48 algorithm are used to evaluate the model's performance. According to the results, the J48 classifier of the proposed model achieves high accuracy comparing to other classifiers. This is because of using the resample filter for the pre-processing phase in the proposed model rather than feature selection technique that used in [9] as illustrated in Table 4.

**Table 4.** Compression of accuracy measures for the Breast Cancer Dataset.

Methodology	Study [9]	Proposed method
With out pre-processing	None	J48: 75.52%, NB: 71.67% SMO: 69.58%
With pre-processing	Missing values were replaced with WEKA pre-processing techniques and feature selection was applied J48: 79.97%, MLP: 75.35% & rough set: 71.36%	Delete records of missing values and Descretization J48: 74.82%, NB: 75.53% SMO: 72.66%
Using the resample filter	None	Applying the resample filter for 7 times J48: 98.20%, NB: 76.61% SMO: 95.32%

## 5.2 Experiment Using the WBC Dataset

Same experiments were applied with the WBC dataset. With respect to applying preprocessing techniques all algorithms present higher classification accuracy, the difference lies in the fact that using the resample filter several times improves the classification accuracy. SMO classifier achieve 99.56% efficiency compared to 99.12% of the Naïve Bayes and 99.24% of the J48. Results are illustrated in Table 5.

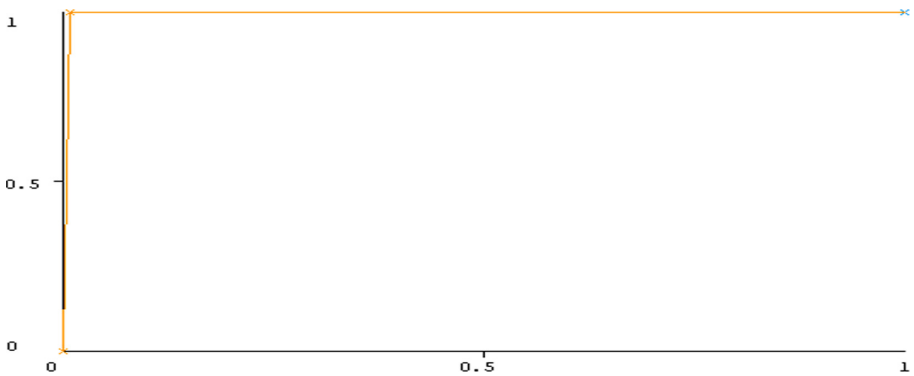
**Table 5.** Performance of the classifiers in WBC dataset.

Experiments steps	Classifier accuracy		
	J48	NB	SMO
Original Without preprocessing	94.56%	95.99%	96.99%
After removing missing values & discretization	95.91%	97.37%	96.78%
After applying resample filter (first time)	95.91%	97.51%	98.97%
Applying resample filter (second time)	97.95%	98.10%	99.41%
Applying resample filter (third time)	98.68%	98.10%	99.12%
Applying resample filter (fourth time)	99.24%	99.12%	99.56%

In the WBC dataset, SMO superior than others with 99.56%. Accuracy measures for SMO classifier is shown in Table 6 and Roc curve of SMO is shown in Fig. 3.

**Table 6.** Accuracy measures for SMO in WBC Dataset.

TP	FP	Precision	Recall	Roc curve	Std	Class
0.996	0.004	0.998	0.996	0.996	0.2220	Benign
0.996	0.004	0.992	0.996	0.996		Malignant



**Fig. 3.** SMO ROC curve in WBC Dataset.

In terms of the WBC dataset, our proposed method is compared with two studies [6, 10]. Results shows that the performance of SMO classifier is better since our model employs pre-processing, and resampling approaches. Thus, utilizing pre-processing, and resampling techniques play an important role in increasing the SMO accuracy comparable to the other techniques in [6, 10]. Details are shown below in Table 7.

**Table 7.** Compression of accuracy measures for the WBC Dataset.

Methodology	Study [6]	Study [10]	Proposed method
Without pre-processing	C4.5: 95% NB: 95.9% SVM: 97.3%	SMO: 96.19%, IBK: 95.90%, BF Tree: 95.46%	J48: 94.56%, NB: 95.99% SMO: 96.99%
With pre-processing	None	None	Delete records of missing values and Descrretization J48: 95.91%, NB: 97.37% and SMO: 96.78%
Using the resample filter	None	None	Applying the resample filter for 4 times J48: 99.24%, NB: 99.12%, SMO: 99.56%



## 6 Conclusion

Breast cancer is considered to be one of the significant causes of death in women. Early detection of breast cancer plays an essential role to save women's life. Breast cancer detection can be done with the help of modern machine learning algorithms. In this paper, we focus on how to deal with imbalanced data that have missing values using resampling techniques to enhance the classification accuracy of detecting breast cancer. In our work, three classifiers algorithms J48, NB, and SMO applied on two different breast cancer datasets. Results show that using the resample filter in the preprocessing phase enhances the classifier's performance. In the future, the same experiments will apply to different classifiers and different datasets.

## References

1. U.S. Cancer Statistics Working Group. United States Cancer Statistics: 1999–2008 Incidence and Mortality Web-based Report. Atlanta (GA): Department of Health and Human Services, Centers for Disease Control
2. [http://www.breastcancer.org/symptoms/understand\\_bc/statistics](http://www.breastcancer.org/symptoms/understand_bc/statistics)
3. Silva, J., Lezama, O.B.P., Varela, N., Borrero, L.A.: Integration of data mining classification techniques and ensemble learning for predicting the type of breast cancer recurrence. In: Miani, R., Camargos, L., Zarpelão, B., Rosas, E., Pasquini, R. (eds.) GPC 2019. LNCS, vol. 11484, pp. 18–30. Springer, Cham (2019). [https://doi.org/10.1007/978-3-030-19223-5\\_2](https://doi.org/10.1007/978-3-030-19223-5_2)
4. Ojha U., Goel, S.: A study on prediction of breast cancer recurrence using data mining techniques. In: 7th International Conference on Cloud Computing, Data Science & Engineering-Confluence, IEEE, pp. 527–530, 2017
5. Pritom, A.I., Munshi, M.A.R., Sabab, S.A., Shihab, S.: Predicting breast cancer recurrence using effective classification and feature selection technique. In: 19th International Conference on Computer and Information Technology (ICCIT), pp. 310–314. IEEE (2016)
6. Asri, H., Mousannif, H., Al, M.H., Noel, T.: Using machine learning algorithms for breast cancer risk prediction and diagnosis. *Procedia Comput. Sci.* **83**, 1064–1069 (2016)
7. Hazra, A., Mandal, S.K., Gupta, A.: Study and analysis of breast cancer cell detection using Naïve Bayes, SVM and Ensemble Algorithms. *Int. J. Comput. Appl.* **145**, 0975–8887 (2016)
8. Rodrigues, B.L.: Analysis of the Wisconsin Breast Cancer dataset and machine learning for breast cancer detection. In: Proceedings of XI Workshop de Visão Computacional, pp. 15–19 (2015)
9. Saabith, A.L.S., Sundararajan, E., Bakar, A.A.: Comparative study on different classification techniques for breast cancer dataset. *Int. J. Comput. Sc. Mob. Comput.* **3**(10), 185–191 (2014)
10. Chaurasia, V., Pal, S.: A novel approach for breast cancer detection using data mining techniques. *Int. J. Innovative Res. Comput. Commun. Eng.* **2** (2017). (An ISO 3297: 2007 Certified Organization)
11. Salama G.I., Abdelhalim, M.B., Zeid, M.A.E.: Experimental comparison of classifiers for breast cancer diagnosis. In: 2012 Seventh International Conference on Computer Engineering & Systems (ICCES), pp. 180–185. IEEE (2012)
12. Lavanya, D., Rani, D.K.U.: Analysis of feature selection with classification: breast cancer datasets. *Indian J. Comput. Sci. Eng. (IJCSE)*, pp. 756–763 (2011)
13. Breast Cancer Wisconsin Dataset. Available at: UCI Machine Learning Repository
14. Dataset Description. Available at: UCI Machine Learning Repository
15. Han, J., Pei, J., Kamber, M.: *Data Mining: Concepts and Techniques*. Elsevier, New York (2011)

16. Quinlan, R.C.: 4.5: Programs for Machine Learning. Morgan Kaufmann Publishers Inc., San Francisco (1993)
17. Quinlan, J.R.: Simplifying decision trees. *Int. J. Man-Mach. Stud.* **27**, 221–234 (1987)
18. Piatt, J.: Fast training of support vector machines using sequential minimal optimization. In: *Advances in Kernel Methods-Support Vector Learning* (1998)
19. Darrab, S., Ergenc, B., Vertical pattern mining algorithm for multiple support thresholds. In: *International Conference on Knowledge Based and Intelligent Information and Engineering (KES)*, *Procedia Computer Science*, vol. 112, pp. 417–426 (2017)
20. Darrab, S., Ergenc, B.: Frequent pattern mining under multiple support thresholds, the International Conference on Applied Computer Science (ACS). *Wseas Transactions on Computer Research*, pp. 1–10 (2016)
21. Alghodhaifi, H., Alghodhaifi, A., Alghodhaifi, M.: Predicting Invasive Ductal Carcinoma in breast histology images using Convolutional Neural Network. In: *2019 IEEE National Aerospace and Electronics Conference (NAECON)*, pp. 374–378 (2019)