

# Effective breast Tumor Classification using K-Strongest Strength With Local Outlier Factor Algorithm

Sannasi Chakravarthy S R<sup>1</sup>, Harikumar Rajaguru<sup>2</sup>

<sup>1</sup> Department Of Electronics And Communication Engineering, Bannari Amman Institute Of Technology, Sathyamangalam – 638401, India.

<sup>2</sup> Department Of Electronics And Communication Engineering, Bannari Amman Institute Of Technology, Sathyamangalam – 638401, India.

Email: <sup>1</sup>elektroniqz@Gmail.Com, harikumarrajaguru@Gmail.Com<sup>2</sup>

**Abstract:** *Eventhough People In This World Are Well-Educated And Sophisticated, Cancer Is Still A Deadly Disease Throughout The World. Amidst This, Breast Cancer Is A Major Cause Of Mortality Among Women. This Shows That There Is Always A Need For The Earlier Tumor Detection Of Breast Cancer. The Paper Utilizes The K-Strongest Strength (Kss) Algorithm for Breast Cancer Detection. The Employed Kss Algorithm Is Influenced by The Law Of Universal Gravitation Analogy And Is Characterized Similarly To The Standard K-Nearest Neighbor (Knn) Algorithm. The Algorithm Is Evaluated Using The Dataset Of Breast Cancer Wisconsin Classification (WDBC) Data. This Input Data Is Preprocessed And Checked For Any Outliers Followed By Their Removal. Thereafter, The Preprocessed Data Is Applied With The Kss Algorithm For Getting A Better Result Of 97.08% Accuracy. The Obtained Results Are Then Compared With The Standard Benchmark Algorithms Such As Knn And Multi-Layer Perceptron Algorithms For Checking The Robustness Of The Kss-LOF Classifier.*

**Keywords:** *Breast Tumor, Breast Cancer, Knn, Kss, Wdbc, Malignant, Outliers.*

## 1. INTRODUCTION

Cancer, Otherwise Referred To As Malignancy That Indicates The Abnormal Growth Of Human Cells. Commonly, More Number Of Tumor Diseases Are Identified In Humans At Different Body Parts Such As Lung Tumor, Breast Tumor, Prostate Tumor, Colon Tumor, Lymphoma, And Skin Tumors [1]. In All These Tumor Diseases, The Symptoms Might Vary. Out Of All The Abovesaid Tumor Diseases, Breast Cancer Is The Salient One That Is Responsible For Causing Higher Mortality Among Women. That Is Why Breast Cancer Is The Second Most Cancer Type Next To Lung Cancer Globally [2]. All The Cancer Diseases Are Diagnosed And Named With Reference To The Cell Organ It Originates. Accordingly, Breast Tumors Start From The Breast Cells Of Women. Additionally, It Is Found In The Literature That The Breast Cancer Incidence Is More For Women Than Men [3].

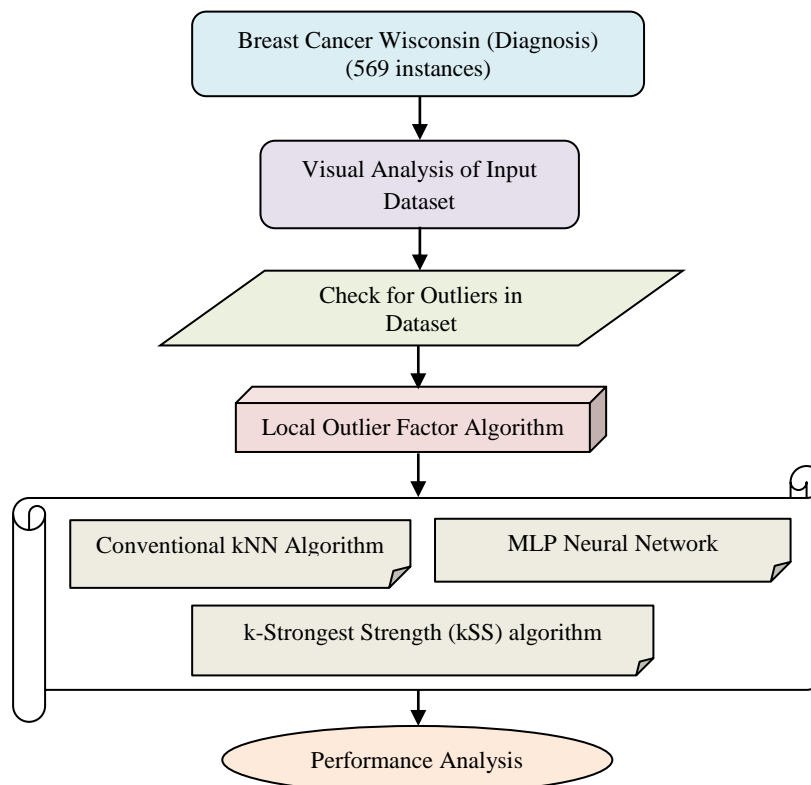


Fig. 1. Proposed Approach For Breast Cancer Classification

From The Above Discussion, It Is Clear That The Impact Of Breast Cancer On The Lifespan Of Women Is More And So, A Robust Methodology Is Always Demanded For Breast Cancer Classification. In This Way, The Researchers Can Help Society To Improve The Lifespan Of Victims, And Thereby The Aim Is To Depreciate The Mortality Rate Of Breast Tumors. Several Research People Are Proposing Newer Methodologies On This Problem For Supporting Human Lives [4]. In Breast Tumor Malignancy, A Lump Is Formed In The Breast Cells And Is Considered As The Major Symptom [5]. The Ducts And Lobules Of The Breast Part Are Affected More At The Initial Stage Of This Cancer Type, But The Above Symptoms Are Not Enough To Confirm The Disease Due To Their Harder Visibility [6]. The Work Evaluates Its Effectiveness In The Breast Cancer Wisconsin Classification (WDBC) Dataset Available At The UCI Repository [7]. In This Dataset, The Characteristics Of The Breast Were Abstracted And Examined Over The Fine Needle Aspect Of The Biopsy Method. The Work Followed For The Breast Cancer Classification In This Paper Is Illustrated In Figure 1.

As Given In Fig. 1, The Work Puts To Good Use Of The WDBC Dataset Comprises 569 Instances Of Breast Study. Then The Visual Analysis Is Performed In The Input Dataset For Better Understanding. From This, The Outliers Are Detected And They Are Analyzed Using The Local Outlier Factor (LOF) Algorithm. Then The Outliers Are Removed And The Classification Phase Is Done With The Final Analysis Of The Performance Of The Kss Algorithm. All The Above Works Are Implemented Using Python 3.6 In Windows 10 Environment Having 8 GB RAM And 2 TB Mass Storage.

## 2. INPUT DATASET AND DATA PREPROCESSING

### 2.1 Dataset Used

The Breast Cancer Wisconsin (Classification) Dataset Is Utilized In This Paper For Finding

The Effectiveness Of The Classification Algorithms. The WDBC Dataset Is A Standard As Well As Publicly Available One At The UCI Repository. The Data Corpus Contains A Total Of 569 Readings With The Numerical And Predictive Data Attributes [7]. Among All, Few Attributes Are Texture, Radius, Area, Perimeter, Compactness, Smoothness, Symmetry, Fractal Dimension, Concavity, And Concave Points. In This, The Local Variations Procured During The Biopsy Procedure Using Fine-Needle Aspect Are Considered As The Smoothness Feature Whereas Another Feature Named Compactness Is Described As,

$$\text{Compactness} = \frac{\text{perimeter}^2}{\text{area}-1} \quad (1)$$

Another Important Feature, Concave Represents The Severity Nature Of Concave Point Calculations On Its Contours. In Addition, The Knowledge On Concave On Its Own Contour Is Considered As A Concave Point Attribute. Consequently, The WDBC Database Was Presented Publicly As Elaborated In [7]. The Beauty Of The Dataset Lies In The Absence Of Missing Values Which Made This WDBC The Most Prevalent Dataset Among Research People. As A Whole, The Database Of WDBC Comprises Two Different Severity Targets Namely Benign (B) Cases And Malignant (M) Targets. Thus, The Paper Intends For The Binary Classification Problem Of Breast Cancer. In Short, The Dataset Has 357 B Cases And 212 Malignant Cases.

## 2.2 Graphical Analysis Of Input Features

Fig.2 illustrates The Visual Analysis Of The WDBC Dataset Using A Boxplot where The X-Axis Denotes The Feature Attributes And The Y-Axis Represents The Feature Values. Here, Both The Axis Are Plotted Against The Two Output Targets, B And M, And B Is Represented As 0 Whereas M Is Represented As 1. As Denoted In Fig. 2, The Relationship Among The Feature Attributes Is Revealed By The Boxplot. That Is, The Plot Reveals That The Input Features Are More Non-Linear, And Also The Plot Shows The Presence Of Outliers In The Input. Additionally, Fig. 2 Is Plotted Using The Seaborn Library With StandardScaler For Easier Data Visualization.

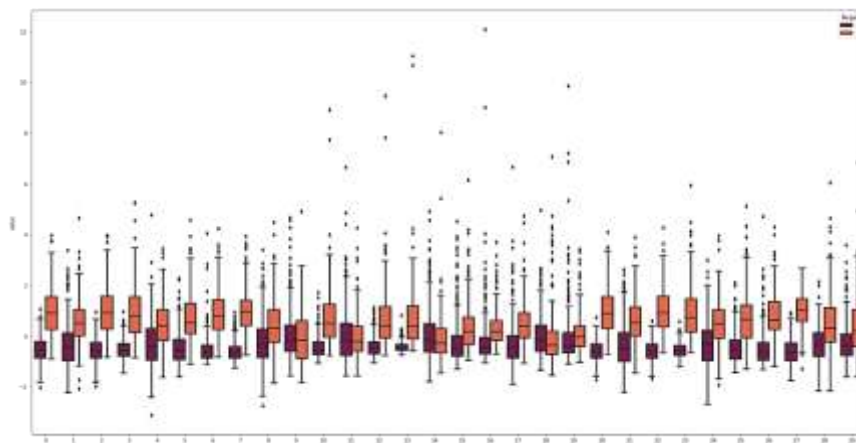


Fig. 2.Box Plot Visualization Of WDBC Attributes

## 2.3 Detecting Outliers Using Local Outlier Factor (LOF) Algorithm

The LOF Approach Is An Unsupervised Method That Follows An Anomaly Detection Procedure That Is Based On The Computation Of Deviation In The Local Density Of The Input Data Over Its Neighbors [8]. For This, The Work Uses Scikit-Learn Package For Its Implementation. Herein, The Paper Typically Sets The Parameter, Number Of Neighbors In Two Ways; The First One Is Higher Than The Minimum Sample Features A Cluster Had,

And Lower Than The Maximum Sample Features That Are Prevalent To Be Local Outliers. In The Above First Method, The Other Sample Features Might Be A Local Outlier When Compared With This Cluster. Accordingly, The Number Of Neighbor Parameter Is Optimally Set As 20 For Our Work.

In This Way, The Outliers In The Input Data Are Determined And Plotted In Fig. 3. As Shown, A Count Plot Is Plotted That Reveals That 58 Outliers And 511 Inliers In The Input Dataset. All These 58 Outliers Score Is Calculated, Sorted And Some Of Them Are Removed Based On The Threshold Of Negative 2.5 Value. This Is Portrayed In The Plot Of Fig. 4. In This Plot, The Outliers Are Represented As Larger Black Dots, The Original Data Points Are Denoted As Smaller Red Dots, And The Outliers Are Encircled Based On The Outlier Scores. Also, The Outlier With More Threshold Are Encircled With Higher Diameter Circle. Here, Based On The Calculated Outlier Scores, The Data Points Are Encircled And Will Be Dropped For The Further Classification Stage So That A Better Result Will Be Obtained.

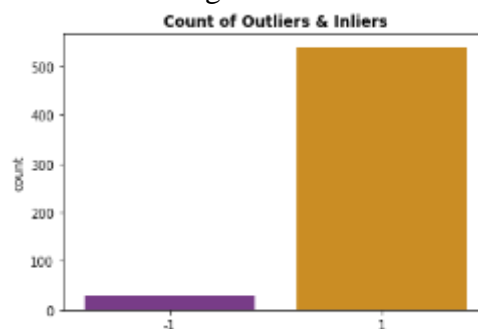


Fig. 3. Outliers And Inliers Count In The Input Data using LOF Algorithm

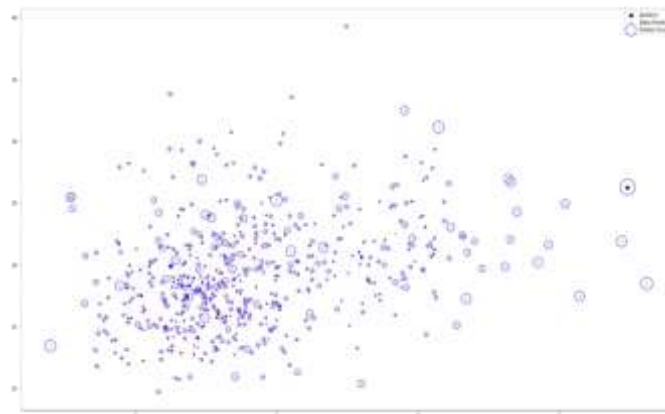


Fig. 4. Outliers With A Threshold Score Of Negative 2.5 Value

### 3. CLASSIFICATION ALGORITHMS

#### 3.1 K-Nearest Neighbor (KNN) Algorithm

The K-Nearest Neighbor (Knn) Model Is A Simple And Most Commonly Used Supervised Algorithm That Need Not Have Any Assumption On The Inputs [9]. This Makes The Algorithm As A Non-Parametric One And Is Thus Referred To As A Lazy Model. The Lazy Learning Algorithm Does Not Involve Immediate Learning From The Training Inputs But It Saves The Data Points And Performs Necessary Action On The Data Points During The Classification Phase [9]. The Steps Of The Knn Algorithm Are Summarized Below [9]:

- #1: Choose The Value Of  $K$  I.E. Number Of Neighbors
- #2: Euclidean Distance Is Computed For The Above  $K$ -Value
- #3: Select The  $K$  Closest Neighbors Based On The Determined Euclidean Distance.
- #4: In This Selection, Count The Data Points That Belong To A Particular Class (B Or M).
- #5: Classify The New Data To A Particular Class For Which The Value Of  $K$  Is Maximum.

### 3.2 Multi-Layer Perceptron (MLP) Algorithm

The Artificial Neural Network (ANN) Models are the ones that are used popularly in several fields of research. Unlike other classifiers, the MLP is an ANN model used for the process of classification tasks [10]. Since it is a feed-forward ANN algorithm, the MLP can map the input set onto its respective targets. An MLP architecture comprises of a multi-layer network where each layer is well-connected with others [10]. The nodes present in each layer are neurons that utilize a non-linear activation map for performing the classification. This ANN type comprises one input, one hidden, and one output layer in which it performs the classification done in the hidden layer.

### 3.3 The $K$ -Strongest Strengths (Kss) Algorithm

The network behind the working of the Kss algorithm depends on the hypothesis in which a better performance in terms of classification will be attained if data points in the training inputs are characterized with a mass, that can depend on their significance and this mimics the physical attracting forces possessed by bigger ones [11]. For the implementation part of the Kss algorithm, the pseudo-code can be summarized as below [11]:

**Input:** Test The Objects For Classifying ( $x_{in}$ ) The Input Training Set Of  $X = \{x_1, x_2 \dots x_n\}$

**Output:** For An Output Target Class,  $x_{in}(f(x_{in}))$

#1: Calculate Mass For Each Training Input Vectors  $x_{in} \in X$ .

#2: **For** all  $i = 1$  **to**  $n$  **do**

#3: Allocate A Mass For This Training Set  $x_{in}$

#4: Here, *Topology-Based* Mass Function Can Be Used

#5: **End For**

#6: Determine The Force Of Attraction Between  $x_{in}$  For Every  $X$  Input

#7: **For** all  $i = 1$  **to**  $n$  **do**

#8: **Strength**( $x_{in}, X$ ) =  $\frac{m(x_{in})}{dist(x_{in}, X)^2}$

#9: **End For**

#10: Determine The Sub-Sets  $s_{x_{in}} \subseteq X = \{x_1, x_2 \dots x_n\}$  Of  $k$  Training Targets That Exert The Dominancy On  $x_{in}$ .

#11: Compute  $f(x_{in}): f(x_{in}) \leftarrow \text{argmax}_{c \in \text{Output Target}} \sum_{i=1}^k \delta(c, x_{in})$

#12: Return  $f(x_{in})$

In the above procedure of the Kss classifier, deprived of computing the distance as in the case of the Knn algorithm for finding the neighborhood, an attraction force can be used [11] and determined using the above steps from 7 to 9. This can be done using the law of universal gravitation as depicted from Newton's discovery in the year 1687.

## 4. RESULTS AND DISCUSSION

The Input Dataset Taken For The Problem Is WDBC Data, Splitted Randomly Using The Standard 70:30 Ratio Of Training And Testing Inputs. The Abovesaid Classification Methodologies Are Implemented And Its Corresponding Training And Testing Are Done Through These Input Sets. Also, A Five-Fold Cross-Validation Approach Is Used For Obtaining The Results. For The Implementation Part, Python 3.6 Is Used As An IDE Environment. The Results Are Obtained After The Implementation And Their Performance Is Analyzed Through The Use Of Benchmark Metrics Such As Sensitivity (Sen), Accuracy (Acc), Specificity (Spe), F1 Score, Precision (Pre), And Matthews Correlation Coefficient (MCC) [12][13][14]. These Performance Measures Are Derived From The Contingency Table Concept That Is Widely Used In Both Binary As Well As Multiclass Classification Problems.

Table 1. Confusion Matrix Obtained For The Test Dataset

Classification tool	Confusion Matrix			
	TP	FN	FP	TN
K-Nearest Neighbor (Knn) Algorithm	56	8	10	97
Multi-Layer Perceptron (MLP) Network	57	7	7	100
K-Strongest Strengths (Kss) Algorithm	62	2	3	104

Table 1 illustrated Above Portrays The Obtained Confusion Matrix Values For different Classification Algorithms used For The Classification Task Of Breast Cancer. As Given In Table 1, More Amount Of Pseudoclassification Results Is Obtained For The Conventional Knn algorithm applied With The Local Outlier Factor Approach, Similarly, The More Amount Of True classifications Is Obtained For The Kss algorithm applied With The Local Outlier Factor Approach. The Values As Depicted In The Confusion Matrix Of Table 1 are Attained and Evaluated for Both Benign Cases And Malignant Classes.

Table 2 Illustrates The Performance Analysis And Comparison Of Distinct Classification Algorithms Through The Use Of The Obtained Confusion Matrix Values As Mentioned In Table 1. As In Table 2, Six Different Performance Metrics Are Utilized For The Comparative Analysis Of The Performance Of The Adopted Classification Tools For Our Breast Cancer Problem.

Table-2. Comparative analysis Of Different Classification frameworks

Classification Tools	Performance Metrics					
	Sen (%)	Spe (%)	Acc (%)	Pre (%)	F1 Score (%)	MCC (%)
K-Nearest Neighbor (Knn) Algorithm	87.5	90.65	89.47	84.85	86.15	77.69
Multi-Layer Perceptron (MLP) Network	89.06	93.46	91.81	89.06	89.06	82.52
K-Strongest Strengths (Kss) Algorithm	96.88	97.20	97.08	95.38	96.12	93.78

The Comparison Of The Performance Of The Algorithms Equipped For The Classification Process Is Summarized And Plotted Graphically in Fig.5. The Accuracy Of Classification For The Kss Algorithm With The Local Outlier Factor Approach Is Obtained as Higher Among Others As Given In Fig. 5. Here, The Maximum Classification Accuracy Of 97.08% Is Obtained For This Kss algorithm Together With The Local Outlier Factor Approach.

Although the Traditional kNN Algorithm being A Lazy Model, The Algorithm Provides An Accuracy Performance Of 89.47%. This Is Possible Due To The Local Outlier Factor Approach utilized For The Outlier Detection And Removal In The Dataset. On The Other Hand, The Standard MLP Neural Network Architecture Produces A Performance In Which Its Classification Strategy Lies In The Middle Of Both Knn And Kss Algorithms. The MLP Network Provides Maximum Accuracy Of 91.81% After Applying The Local Outlier Factor Approach.

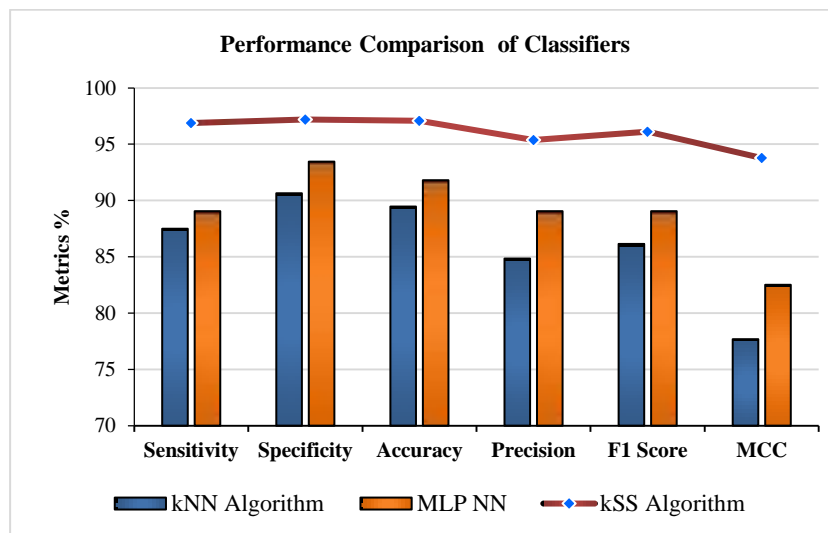


Fig. 5. Graphical plot Of Performance Of Distinct classification Algorithms

As Shown In Fig. 5 And Table 2, The Classification Performance Of The Kss With The LOF Algorithm Is Higher Than Other Classifiers Used. As In Fig. 5, The Kss Algorithm Approach Of The Classification Framework Provides A Superior Classification Over The Knn And MLP Frameworks. This Reveals That The Kss Algorithm Is Good In Discriminating The Severities Of Benign And Malignant For Breast Cancer Classification. Also, Reveals That The Superior Performance Of The Kss Algorithm Is Validated Through The MCC Obtained As 93.78% With The Precision And F1 Score Of 95.38% And 96.12% Respectively.

## 5. CONCLUSION

The Work Involves The Design Of A Robust Computer-Aided Tool For Classifying The Severities In Breast Cancer. The Aim Is To Classify The Input Dataset Vectors Into Either Benign Class Or Malignant Target. Herein, The WDBC Dataset From The UCI Repository Is Employed For The Dissemination Of The Work. The Dataset Is Visually Analyzed For Its Non-Linearity And The Presence Of The Outliers Detected In The Input Data Using The Local Outlier Factor Algorithm. Now The Outliers With A Maximum Outlier Score Are Removed And Then The Cleaned Data With No Outliers Are Moved On To The Classification Stage. Here, Three Classification Algorithms Are Employed Namely, The Knn Algorithm, MLP Neural Network, And Kss Algorithm For The Intention To Classify The Inputs Into Their Respective Severities, Benign Or Malignant. Thus, The Kss Algorithm Based On The Law Of Universal Gravitational Strategy And In Par With The Similar Operation Of The Knn Algorithm Together With The Idea Of Outliers Detection And Removal, The Work Attains A Maximum Classification With 97.08% Accuracy With An MCC Validation Score Of 93.78%. The Future Work Will Be In Applying The Kss Algorithm With The LOF Approach For Other Input Datasets With Different Pre-Processing



Strategies.

## 6. REFERENCES

- [1] Sung, H., Ferlay, J., Siegel, R.L., Laversanne, M., Soerjomataram, I., Jemal, A. And Bray, F., 2021. Global Cancer Statistics 2020: GLOBOCAN Estimates Of Incidence And Mortality Worldwide For 36 Cancers In 185 Countries. CA: A Cancer Journal For Clinicians, 71(3), Pp.209-249.
- [2] Dibden, A., Offman, J., Duffy, S.W. And Gabe, R., 2020. Worldwide Review And Meta-Analysis Of Cohort Studies Measuring The Effect Of Mammography Screening Programmes On Incidence-Based Breast Cancer Mortality. Cancers, 12(4), P.976.
- [3] Cao, W., Chen, H.D., Yu, Y.W., Li, N. And Chen, W.Q., 2021. Changing Profiles Of Cancer Burden Worldwide And In China: A Secondary Analysis Of The Global Cancer Statistics 2020. Chinese Medical Journal, 134(7), P.783.
- [4] Sannasi Chakravarthy, S.R. And Rajaguru, H., 2021. A Novel Improved Crow-Search Algorithm To Classify The Severity In Digital Mammograms. International Journal Of Imaging Systems And Technology, 31(2), Pp.921-954.
- [5] Abirami, C., Harikumar, R. And Chakravarthy, S.S., 2016, March. Performance Analysis And Detection Of Micro Calcification In Digital Mammograms Using Wavelet Features. In 2016 International Conference On Wireless Communications, Signal Processing And Networking (Wispnet) (Pp. 2327-2331). IEEE.
- [6] Sannasi Chakravarthy, S.R. And Rajaguru, H., 2020. Detection And Classification Of Microcalcification From Digital Mammograms With Firefly Algorithm, Extreme Learning Machine And Non-Linear Regression Models: A Comparison. International Journal Of Imaging Systems And Technology, 30(1), Pp.126-146.
- [7] Blake, C.L. And Merz, C.J., 1998. UCI Repository Of Machine Learning Databases. University Of California, Irvine, Dept. Of Information And Computer Sciences.
- [8] Cheng, Z., Zou, C. And Dong, J., 2019, September. Outlier Detection Using Isolation Forest And Local Outlier Factor. In Proceedings Of The Conference On Research In Adaptive And Convergent Systems (Pp. 161-168).
- [9] Rajaguru, H., 2019. Analysis Of Decision Tree And K-Nearest Neighbor Algorithm In The Classification Of Breast Cancer. Asian Pacific Journal Of Cancer Prevention: APJCP, 20(12), P.3777.
- [10] Sannasi Chakravarthy, S.R. And Rajaguru, H., 2021. A Hybrid Classification Framework For The Effective Classification Of Breast Tumors. Journal Of Computational And Theoretical Nanoscience, 18(3), Pp.959-965.
- [11] Demidova, L., Nikulchev, E. And Sokolova, Y., 2016. The Svm Classifier Based On The Modified Particle Swarm Optimization. Arxiv Preprint Arxiv:1603.08296.
- [12] Hossin, M. And Sulaiman, M.N., 2015. A Review On Evaluation Metrics For Data Classification Evaluations. International Journal Of Data Mining & Knowledge Management Process, 5(2), P.1.
- [13] K. Yasoda, R. Ponnagall, K. Bhuvaneshwari, And K. Venkatachalam, "Automatic Detection And Classification Of EEG Artifacts Using Fuzzy Kernel SVM And Wavelet ICA (WICA)," *Soft Computing*, Vol. 24, No. 21, Pp. 16011-16019, 2020.
- [14] C. Viji, N. Rajkumar, S. Suganthi, K. Venkatachalam, And S. Pandiyan, "An Improved Approach For Automatic Spine Canal Segmentation Using Probabilistic Boosting Tree (PBT) With Fuzzy Support Vector Machine," *Journal Of Ambient Intelligence And Humanized Computing*, Pp. 1-10, 2020.