

Effective classification framework For Breast Tumors using Optimized Multi-Kernel Svm With Controlled Skewness

Harikumar Rajaguru¹, Sannasi Chakravarthy S R²

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

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

harikumarrajaguru@gmail.com¹, elektronikz@gmail.com²

Abstract: *The World Is So Advanced And Sophisticated Enough Nowadays. But, Cancer Still Remains As A Deadly Disease In Many Parts Of The World For All Living Organisms. The Incidence Rate Of Cancer Among Humans Globally Is Increasing Steadily Day By Day. Among All Forms Of Tumors, Breast Cancer Is A Type Of Illness That Plays A Major Role In Disturbing the survival rate Of Humans Globally. Thus, There Is A Need To Predict Breast Cancer In Its Earlier Stage. Thus, The Work In This Paper is Aided To Design A Robust Classification Model That Involves The Use Of A Randomized-Parameter Optimized multi-Kernel Support Vector Machine (RPOMK-SVM) Classifier. Before The Stage Of Classification, The Paper Analyzes The Nature Of Input In Order To Obtain Promising Results. The Skewness Of The Feature Attributes Is Controlled And Reformed Using Box-Cox Transform. For This Analysis And Evaluation, The Paper employs the Breast Cancer Wisconsin (Diagnostic) Database, Which Is A Standard Public Dataset. The Final Results Are Then Compared Against The Existing Algorithms.*

Keywords: *Breast Cancer, Support Vector Machine, Malignant, Skewness, Benign, Gaussian, Wdbc.*

1. INTRODUCTION

Cancer, otherwise known as malignancy, refers to the abrupt cell growth in the human body. In general, there are more than a hundred different cancer types identified, some of them are lung cancer, breast cancer, colon cancer, prostate cancer, skin cancer, and lymphoma [1]. The symptoms may vary based on the cancer type. In this, breast tumor is the one type of cancer that leads to more mortality among women. Besides, breast cancer remains in second position among all types of cancer, next to lung cancer [2]. Cancers could be named and diagnosed based on the type of cell it originates. In this way, breast cancer starts from the cells of the breast of the human body. In addition, the occurrence frequency of breast cancer is more for women rather than males [3].

As discussed above, breast cancer is the one that heavily affects the survival rate of women. Thus, it is necessary to identify the breast tumor in an earlier way. This will

Increase The Life Span And Will Reduce The Mortality Rate Of Breast Cancer Definitely [4]. Many Researchers Are Working Towards This To Support Human Lives. In Breast Cancer, The Major Symptom Is The Lump Formation In The Cells Of The Breast [5]. During The Initial Stage Of Breast Cancer, The Primary Region Of Ducts And Lobules Are Affected. These Symptoms Are Very Hard To Feel And So The Woman Affected By Breast Cancer Cannot Able To Realize Or Recognize On Its Own [6]. The Paper Utilizes The Standard, As Well As Publicly Available Breast Cancer Wisconsin (Diagnosis) (WDBC) Database. Herein, The Breast Characteristics Were Examined and Abstracted Through A Fine Needle procedure At The Time Of Biopsy. Figure 1 Portrays The Workflow Followed For The Aim Of Effective Classification Of Breast Cancer.

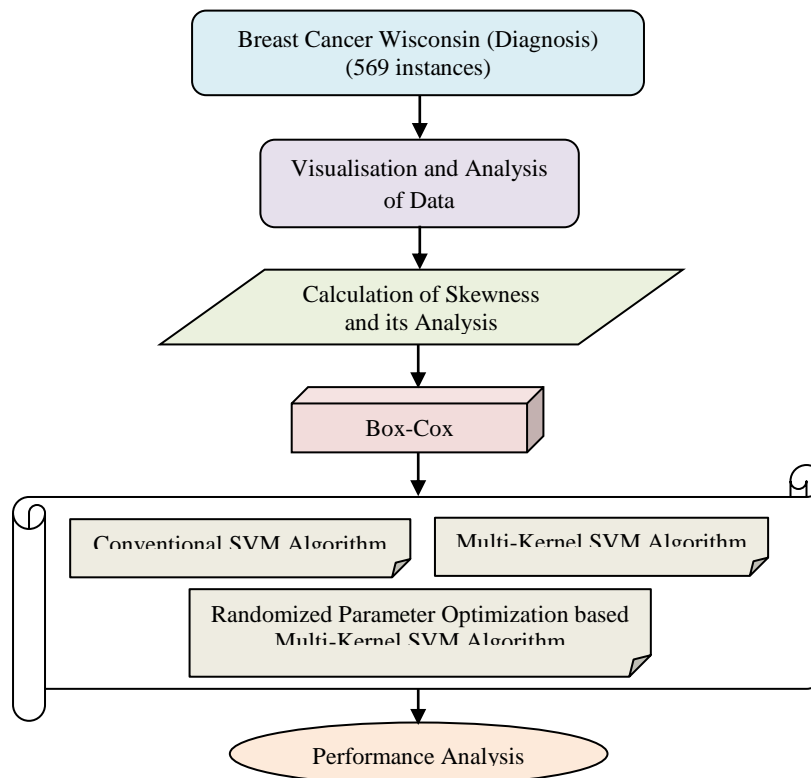


Fig. 1. The Work-Flow Of The Paper

As Portrayed In The Above Figure 1, The Work Makes Use Of Wdbcdata Corpus comprises Of 569 Instances For Its Evaluation. The Input Features Of The WDBC Dataset Are Visualized For Better Analysis. From This, The Skewness Is Calculated And Plotted For The Input Feature Vectors. Then The Box-Cox Transform Is Used For The Correction Of Controlled Skewness Of The Dataset. The Classification Process Is Next Carried Out And Thereby The Results Are Compared And Analyzed.

2 . PREPROCESSING OF DATASET

2.1 Data Visualization and Analysis

As Portrayed In Figure 1, The Breast Cancer Wisconsin (Diagnosis) Data-Corpus Is Employed For The Breast Cancer Classification. This Publicly available Data Is Widely Known In The Name Of WDBC Dataset. And It comprises Of Total Instances Of 569 Readings Together With Predictive (30) And Numerical Attributes [7]. The Sample

Attributes That Are Included In The wdbc Areradius, Texture, Area, Compactness, Perimeter, Symmetry, Smoothness, Concavity, Concave Points, And Fractal Dimension. Herein, During The Biopsy Testconducted With A Fine-Needle Procedure, The Obtained Local Variation In Radius Length Is Said To Be Smoothness Attribute. The Compactness Is Another Attribute That Is Computed As,

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

The Next Concave Attributeillustratesthe Severityof The Calculation Of Concave Portions On Its Own Contour. Also, Themasure Ofthese Concave Detailsthat Are Attained Onits Contour Is Saidto Bea Concave Point Features. In This Way, The WDBC Dataset Was Publicly Introducedas Detailed In [7]. No Missing Values Found In Thedata Set Having 569 Instances Whichmadethis WDBC Data Setmore Popular Among The Cancer Researchers. The Dataset - WDBC Consists Of Two Output Or Severity Targets That Are Denoted As Benign (B) Target And Malignant (M) Target Class. Hence, Theworkaims Toperform Binary Classification After Analyzing The Dataset. The Graphical Distribution Of Output Classes (Severity) Of The WDBC Dataset Is Plotted In Figure 2.

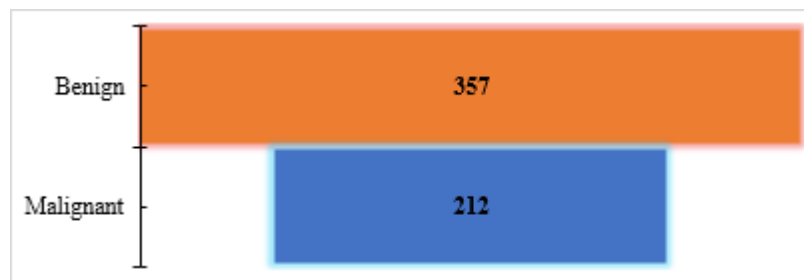


Fig. 2.The Distribution Of Output Severity Present In WDBC Dataset

Figure 3 Portrays The Pair-Plot Visualization Of WDBC Datasetpertaining Tosample Attributes(Mean Radius, Mean Texture, Mean Perimeter, Mean Area, Mean Smoothness)Plotted Against The Output Classes. As Portrayedin Figure 3, The Scatter Plot Reveals That The Input WDBC Features Are Highly Non-Linear. Also, The Diagonal Plot In Figure 3 Represents The Kernel Density Estimates (KDE) Of The Input WDBC Dataset That Reveals The Features Are Skewed At A Different Range. Thus, The Skewness Is Needed To Be Checked Before The Process Of Classification.

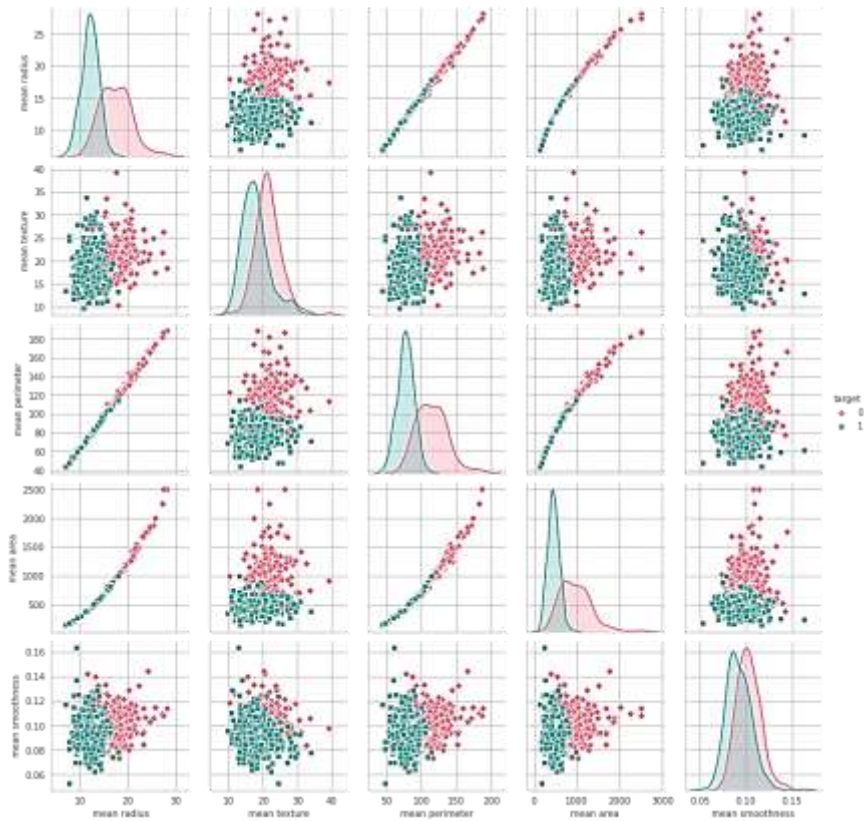


Fig. 3. Visualization Of Sample WDBC Attributes

2.2 Calculation Of Skewness And Its Analysis

In General, Skewness Is A Term That Represents How The Input Data Distracts From The Normal Distribution. In A Normal Distribution, The Samples Are Graphically Denoted As A Bell-Shaped One, Where The Average (Mean) And The Maximum Value In The Dataset (Mode) Are Equal [6]. After Plotting The Input Samples, If The Right-Side Of The Curve Is Found To Be Larger Than Its Left Tail, Then It Can Be Noted That The Input Data Has Positive Skewness I.E. $Mean > Median > Mode$. In Other Cases, That Is, If The Left-Side Of The Curve Is Found To Be Larger Than Its Right Tail, Then It Can Be Noted That The Input Data Has Negative Skewness I.E. $Mode > Median > Mean$. If Both The Left And Right Tails Are Normally Distributed, Then The Data Is Said To Have Symmetric Skewness. The Skewness Is Calculated For All The Attributes Of WDBC Dataset And Is Tabulated In Table 1.

Table 1. Obtained Skewness Value For Different Attributes Of WDBC Dataset

S No	Attributes	Skewness Value	S No	Attributes	Skewness Value
1	Mean Radius	0.942380	16	Compactness Error	1.902221
2	Mean Texture	0.650450	17	Concavity Error	5.110463
3	Mean Perimeter	0.990650	18	Concave Points Error	1.444678
4	Mean Area	1.645732	19	Symmetry Error	2.195133
5	Mean Smoothness	0.456324	20	Fractal Dimension Error	3.923969
6	Mean Compactness	1.190123	21	Worst Radius	1.103115
7	Mean Concavity	1.401180	22	Worst Texture	0.498321

8	Mean Concave Points	1.171180	23	Worst Perimeter	1.128164
9	Mean Symmetry	0.725609	24	Worst Area	1.859373
10	Mean Fractal Dimension	1.304489	25	Worst Smoothness	0.415426
11	Radius Error	3.088612	26	Worst Compactness	1.473555
12	Texture Error	1.646444	27	Worst Concavity	1.150237
13	Perimeter Error	3.443615	28	Worst Concave Points	0.492616
14	Area Error	5.447186	29	Worst Symmetry	1.433928
15	Smoothness Error	2.314450	30	Worst Fractal Dimension	1.662579

The Value Of Skewness As Given In Table 1 Tells Us How The Attributes Are Distorted From The Normal Distribution. The Highest Value Of Skewness Is Found For The Attribute 'Area Error' I.E. 5.447186 And The Least Value Of Skewness Is Obtained For The Attribute 'Worst Smoothness' I.E. 0.415426. In All The Machine-Learning Algorithms, Any Value Of Skewness Is Generally Undesirable, Since This Can Results In an Excessively Large Variance In The Estimates. Thus, For Every Classification Problem, It Is Necessary To Decrease Skewness Value To Make The Data As Closer To A Normal Distribution Curve By Employing any Transformation Method. The Paper Adopts A Box-Cox Transform To Reduce The Skewness In The Input Data.

2.3 Box-Cox Transform

In General, Box-Cox Transform Is Used for The Transformation Of Non-Normal (Skewed) Dependent Values In The Input Data Into A Normal (Bell-Shaped) Shape. This Will Influence Or Enhance The Performance Of Any Classification Framework. The Advantage Of Box-Cox Transform Is That It Is A Configurable Data Transformation procedure that Also Supports The Square-Root And Log-Transformation Methods [8]. Moreover, This Transform Can Be Configurable for automatic Evaluation Of A Suite Of Mathematical Transforms And Thus Provides A Best-Fit For The Input Data. The Box-Cox Transform Can Be Defined As [8],

$$y_i^{(\lambda)} = \begin{cases} \frac{y_i^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \text{Log}(y_i) & \text{if } \lambda = 0 \end{cases} \quad (2)$$

Where The Transform Reduces The Skewness Of Input Data Based On The Box-Cox Parameter (λ). Figures 4 And 5 Show The Skewness Of The 'Area Error' Attribute Before And After Applying The Box-Cox Transform.

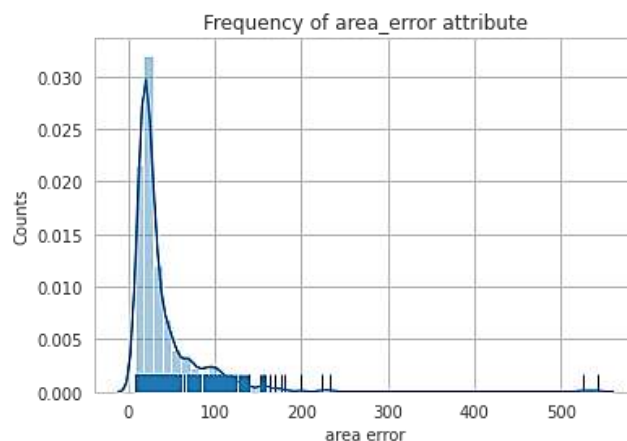


Fig. 4. Skewness Plot Of 'Area Error' Attribute Of WDBC Dataset Before Box-Cox Transform

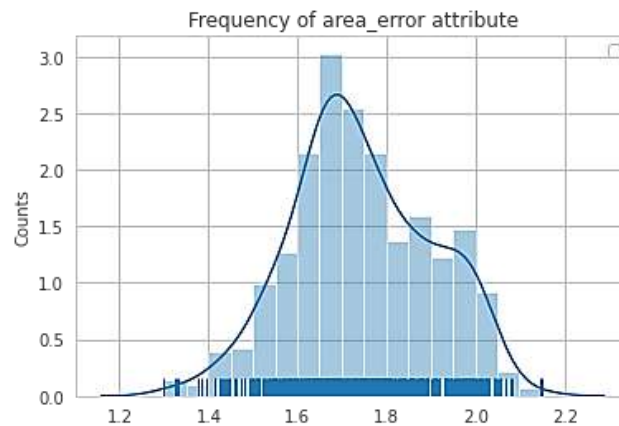


Fig. 5. Skewness Plot Of 'Area Error' Attribute Of WDBC Dataset After Box-Cox Transform

As From Table 1, The Attribute Having A Larger Skewness Is The 'Area Error' Attribute. The Skewness Plot Of The 'Area Error' Plot Is Shown In Figure 4, Where The Skewness Is Not Symmetry. By Using Box-Cox Transform, The Skewness Is Reduced From 5.447186 To 0.058528, And Thus It Resembles The Curve Of Normal Distribution As Shown In Figure 5. In This Way, The Skewness Of All The Attributes Is Reduced To Obtain Better Classification Performance.

3. CLASSIFICATION ALGORITHMS

3.1 Multi-Kernel Support Vector Machine (MK-SVM) Classifier

The Support Vector Machine (SVM) Model Is A Type Of Supervised Learning Algorithm, Used Popularly For The Task Of Classification Problems [9]. The SVM Algorithm Follows A Discriminative Classification Strategy that Is Otherwise Done By Using Separating Hyperplanes. This Implies That The SVM Classifier Provides An Optimal Hyperplane Used For Categorizing The Newer Input Samples. By Using These Hyperplanes, The SVM Can Classify Well Enough For Linear Inputs [10]. But For The Application Of Non-Linear Classification Problems, The SVM Model Makes Use Of Kernel Functions For Categorizing The Newer Input Samples. The MK-SVM Classifier Includes The Combination Of Two Different Kernels And Is Given By,

$$MK - SVM = \frac{1}{2}[(lin(a, b) + rbf(a, b))] \quad (3)$$

Where The *lin* Represents The Use Of Linear Kernel And The *rbf* Represents The Use Of Radial Basis Function In The Multi-Kernel SVM Algorithm.

3.2 Randomized-Parameter Optimized MK-SVM (RPOMK-SVM) Classifier

For Solving Any Type Of Problem, The Exhaustive Grid Search Is Widely Used For Optimizing The Parameters Of SVM. The Parameter Tuning Using An Exhaustive Grid Search Technique Increases The Performance But Decrease The Overall Efficiency Of The System. The Paper Makes Use Of A Randomized Search-Over Parameter Method [11], Where Every Setting Is Obtained As A Distribution Against All The Possible Value Of Parameters. In This Technique, How Parameters Are Needed To Be Sampled Is Carried Out Using A Dictionary As Followed In All Optimization Techniques. For Every Parameter Values, Either A List Of Discrete Choices Or A Distribution Over Possible Values Are Versatile To Optimize The SVM Parameters [12].

Let's Discuss About The Results Of Using The RPOMK-SVM Classifier Along With Box-Cox

Transform To Obtain a Better Classification Of Breast Cancer In The Next Section.

4. RESULTS AND DISCUSSION

The WDBC Data Set is Randomly Divided Into Training (80%) And Testing (20%) Sets. The Above-Mentioned Classification Strategies will Be Implemented And Their Respective Analysis Is Made Using These Training Set And Testing Set. All The Works Discussed Are Done Using Intel (Vpro) Core-I5 Processor, 4 TB Hard-Drive Memory, 8 GB RAM, Installed With Python 3.6 In Windows 7 Operating System. After Obtaining The Classification Results, They Are Analyzed Using Standard Evaluation Metrics - Accuracy (Acc), Sensitivity (Se), Precision (Pr), Specificity (Sp), F1 Score, And Matthews Correlation Coefficient (MCC). The Above Evaluation Metrics will Be Derived From The Concept Of Confusion Matrix, Which Is Used Popularly For Binary And Multiclass Classification Tasks.

Table 2. Obtained Confusion Matrix For Different Algorithms

Classification framework	Confusion Matrix			
	TP	FN	FP	TN
SVM With Box-Cox Transform	168	44	49	308
MK-SVM With Box-Cox Transform	186	26	34	323
RPOMK-SVM With Box-Cox Transform	201	11	16	341

Table 2 shown above shows the values of the Confusion Matrix obtained for distinct Classification Frameworks that are adopted for Breast Cancer Classification Tasks. As in Table 2, the more number of pseudopredictions is attained for the conventional SVM model with Box-Cox Transform and the more number of correct predictions is attained for the RPOMK-SVM with Box-Cox Transform. The Confusion Matrix values as given in Table 2 are obtained and assessed for both severities - Benign and Malignant Outputs.

Table-3. Comparative analysis Of Different Classification frameworks

Classification Framework	Performance Measures (%)					
	Se	Sp	Acc	Pr	F1 Score	MCC
SVM With Box-Cox Transform	79.25	86.27	83.66	77.42	78.32	65.22
MK-SVM With Box-Cox Transform	87.74	90.48	89.46	84.55	86.11	77.65
RPOMK-SVM With Box-Cox Transform	94.81	95.52	95.25	92.63	93.71	89.91

Table 3 portrays the comparison of performance analysis of different Classification Frameworks using the Confusion Matrix elements as shown in Table 2. From Table 3, six distinct evaluation measures are adopted for the performance analysis of the different Classification Frameworks for our Classification Problem.

The Performance Analysis of the Classification Frameworks is calculated and graphically portrayed in Figure 6. As shown in Table 3, the Classification Accuracy is obtained high for the RPOMK-SVM with the Box-Cox Transform framework. In this, the high value of 95.25% Classification Accuracy is yielded for this RPOMK-SVM algorithm together with Box-Cox Transform. Despite the fact that the conventional SVM uses the trick of simple hyperplanes, it provides a Classification Accuracy of 83.66%. This is possible because of the use of Box-Cox Transform employed for the reduction of skewness of

Input Data. The MK-SVM Algorithm With Box-Cox Transform Provides A Better Classification Of 89.46% Of Accuracy. And These Classification Performances Are compared And Analyzed In Table 3 Are Graphically Shown In Figure 6.

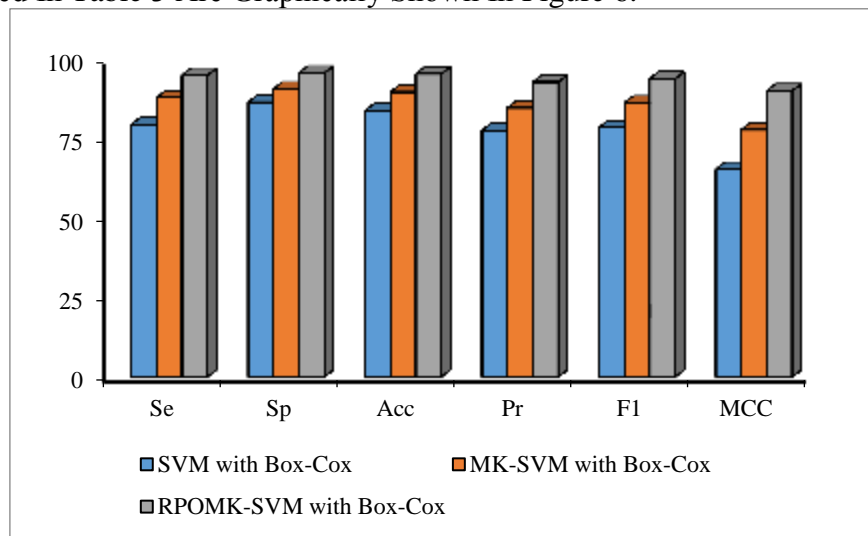


Fig. 6. Graphical Comparison Of Different Classification Frameworks

As Portrayed In Table 3 And Figure 6, The Performance Of The RPOMK-SVM Algorithm Together With Box-Cox Transform Is Significantly Very High Than The Conventional And Multi-Kernel SVM Techniques. As Portrayed In Figure 6, This classification Framework Provides The superior Performance Of Classification That well Differentiate The Benign (B) And Malignant (M) Inputs. That Is, That the RPOMK-SVM Algorithm Together With Box-Cox Transform provides The superior Values Of Sensitivity, F1 Score, Specificity, Precision, MCC, And Accuracy As Compared With Others.

5. CONCLUSION

The Design Of Computer-Aided As Well As Automatic Classification Approach Is Discussed In This Paper For Classifying The benign And Malignant Type Of Severities Pertains To breast Cancer. The Performance Analysis Of The RPOMK-SVM Algorithm Together With The Box-Cox Transform Is Compared And Analyzed With Its Variants Such As Conventional And Multi-Kernel Support Vector Machine Algorithms for The Purpose Of severity Classification Of Benign And Malignant Ones. Herein, The Dataset Used In WDBC Data, And It Is Found That Some Of The Attributes Of The Dataset Are Highly Skewed Up To The Value Of 5.447186. This Skewness Will Affect The Performance Of Any Classification Models And Thus The Box-Cox Transform Is Used To Remove The Skewness Of Input Before Proceeding To The Stage Of Classification. Hence, The randomized-Parameter Optimized Multi-Kernel Support Vector Machine Algorithm Along With The Box-Cox Transform Yields Superior Performance Over Other Techniques. The Future Work Will Be The Use Of Clinical Data For The Designed Classification Framework And For Classifying Other Severities.

6. REFERENCES

- [1] Desantis, C.E., Ma, J., Gaudet, M.M., Newman, L.A., Miller, K.D., Goding Sauer, A., Jemal, A. And Siegel, R.L., 2019. Breast Cancer Statistics, 2019. CA: A Cancer Journal For Clinicians, 69(6), Pp.438-451.

- [2] Northouse, L.L., Templin, T., Mood, D. And Oberst, M., 1998. Couples' Adjustment To Breast Cancer And Benign Breast Disease: A Longitudinal Analysis. *Psycho-Oncology: Journal Of The Psychological, Social And Behavioral Dimensions Of Cancer*, 7(1), Pp.37-48.
- [3] Evans, A., Whelehan, P., Thomson, K., Mclean, D., Brauer, K., Purdie, C., Jordan, L., Baker, L. And Thompson, A., 2010. Quantitative Shear Wave Ultrasound Elastography: Initial Experience In Solid Breast Masses. *Breast Cancer Research*, 12(6), Pp104.
- [4] Chaurasia, V., Pal, S. And Tiwari, B.B., 2018. Prediction Of Benign And Malignant Breast Cancer Using Data Mining Techniques. *Journal Of Algorithms & Computational Technology*, 12(2), Pp.119-126.
- [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] Zhang, Y., Xiong, R., He, H. And Pecht, M.G., 2018. Lithium-Ion Battery Remaining Useful Life Prediction With Box-Cox Transformation And Monte Carlo Simulation. *IEEE Transactions On Industrial Electronics*, 66(2), Pp.1585-1597.
- [9] Dadi, H.S. And Pillutla, G.M., 2016. Improved Face Recognition Rate Using HOG Features And SVM Classifier. *IOSR Journal Of Electronics And Communication Engineering*, 11(4), Pp.34-44.
- [10] Padierna, L.C., Carpio, M., Rojas-Domínguez, A., Puga, H. And Fraire, H., 2018. A Novel Formulation Of Orthogonal Polynomial Kernel Functions For SVM Classifiers: The Gegenbauer Family. *Pattern Recognition*, 84, Pp.211-225..
- [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] Bamakan, S.M.H., Wang, H., Yingjie, T. And Shi, Y., 2016. An Effective Intrusion Detection Framework Based On MCLP/SVM Optimized By Time-Varying Chaos Particle Swarm Optimization. *Neurocomputing*, 199, Pp.90-102.