

Breast Cancer Identification and Classification Using Machine Learning Techniques

Mir Aadil Hussain¹, and Yogesh²

¹M. Tech Scholar, Department of Computer Science & Engineering, RIMT University, Mandi Gobindgarh, Punjab, India
²Assistant Professor, Department of Computer Science & Engineering, RIMT University, Mandi Gobindgarh, Punjab, India

Correspondence should be addressed to Mir Aadil Hussain; raspstest9@gmail.com

Copyright © 2022 Made Mir Aadil Hussain et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

ABSTRACT- The most prevalent disease among women and a major contributor to the rising mortality rate among women is breast cancer. There is a high need for automatic diagnostic systems for early identification of breast cancer since manual breast cancer diagnosis takes a lot of time and there are few systems available. Deep learning and machine learning approaches play a vital role in developing such technologies. We have employed machine learning classification algorithms to distinguish between benign and malignant tumours. These approaches allow the computer to learn from previous data and predict the category of fresh input. Breast cancer is the major cancer in women (43.3 instances per 100,000 women) but the one with the highest mortality rate (14.3 incidents per 100,000 women). For survival, early diagnosis is essential. The issue may be successfully identified, forecasted, and evaluated using machine learning techniques. The eight machine learning methods we compared in this study were the Gaussian Naive Bayes (GNB), kNearest Neighbours (K-NN), Svm Classifiers (SVM), Variational Forest (RF), AdaBoost, Multilayer Perceptron (GB), and datasets from Breast Cancer Badgers. The results of the tests showed that ANN had the better spec. Efficiency (99.28 percent), F1-score (99.99 percent), recall (96.75 percent), accuracy (99.19 percent), and AUC were attained via XGBoost (99,61 percent). Our findings demonstrated that, in the Breast Cancer Wisconsin dataset, ANN is the most successful approach for predicting tumor.

KEYWORDS- Breast cancer, SVM, Random Forest, Gradient boost.

I. INTRODUCTION

One of the foremost causes of death for women is breast cancer. (after lung cancer). In the US, it is anticipated that there will be 246,660 new instances of breast cancer in women in 2016, and that there will be 40,450 deaths of women from the disease. 25 percent of all malignancies in women and around 12 percent of all new cases of cancer are breast cancer.[1] Deep learning and machine learning may have applications in the battle against the disease. Transfer learning and advanced analytics, which have become synonymous with artificial intelligence and data science, have significantly changed the way artificial intelligence is used to make decisions and predict outcomes [2]. In fact, big data has advanced not only the size of data but also generating value from it. Due to its great performance in outcome prediction, lowering medical expenditures, and boosting patients' health while making decisions in real time to save lives, machine learning & deep learning technologies, for example, are being used to medical science problems more and more often. With 43.3 cases per 100,000 women, the most common malignancy in women is breast cancer. Cancer is less common than those of other malignancies. Has a comparatively low mortality rate. It does, however, have the greatest fatality rate of any cancer in women due to the huge number of incidences (12.9 per 100 000). For survival, early diagnosis is essential. Around 70% of cancer fatalities take place in low- and middle-income nations.[3]

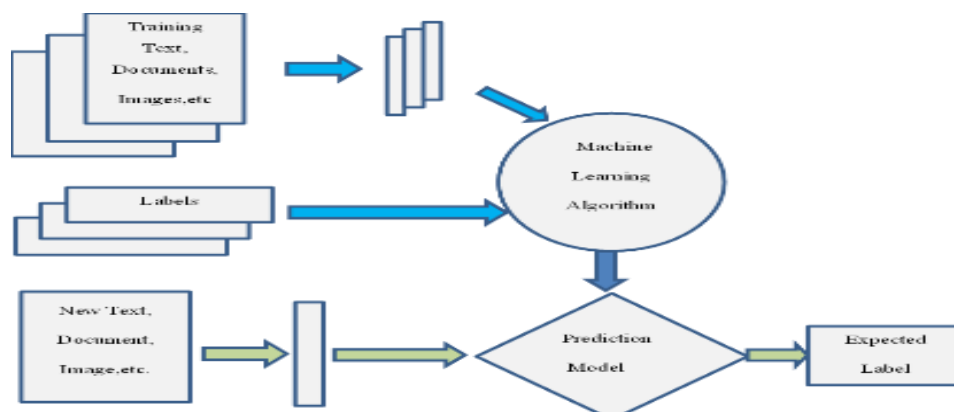


Figure 1: Naive Bayes method

Lack of funding and undeveloped healthcare infrastructures make it challenging for patients to see a doctor. The patient survival rate can be raised by creating strategies for clinical prediction based on early signs and symptoms.; World Health Organization.[5]

Machine learning techniques are more likely to be used to deliver a rapid, automated, and better knowledge of cancer healthcare as the dataset of breast cancer patients grows (Maity, G., and Das, S. 2017).[6] We have the possibility to make an accurate forecast thanks to the enormous datasets that are readily available for detection. What approach will yield the best outcome is the issue, though. Earlier investigations revealed accurate approaches even if earlier study indicated that SVM is, on average, the most accurate method. Amrane, M. et al. (2018) examined Nave Bayes and kNN as two distinct classifiers for the classification of breast cancer. They assessed accuracy using cross-validation techniques. [7]. The outcome demonstrated that Nave Bayes is less accurate than K-NN (97,51 percent). Sharma, S. et al. (2018) evaluated Random Forest or RF, K-NN, and Nave Bayes to predict breast cancer using The Wisconsin Breast Cancer dataset. The analysis finds that KNN has the best performance, at 94.20%. Liu, B. et al. (Liu, B. et al. 2018) looked at a number of techniques, including SVM, AdaBoost, Decision Tree, and Random Forest or RF, in order to spot the both benign and malignant characteristics of malignancy from just a moving photo of a biopsy aspirate of a lumpectomy. The outcomes show that the random forest is the most effective method for projection. [8]

These findings demonstrate that it is possible to compare these methodologies. Using a shared dataset, the optimal approach may be found objectively. There are further techniques to forecast breast cancer that have never been compared. As a result, we AdaBoost, Xgboost (GB), Support vector regression, K-Nearest Neighbors (K-NN), Support Vector Machine (SVM), and Gaussian Naive Bayes were among the machine learning algorithms. whose efficacy was evaluated (GNB).

A. Machine Learning

By giving computers data in the form of instances of real-world occurrences, machine learning enables computers to learn and behave like people do, improving their behavior over time in an independent manner. Our arsenal of machine learning techniques includes clustering, artificial neural networks, decision trees, and Bayes networks, among others. The preferred method relies on the problem that has to be solved and the amount of information available.[9]. Structural machine learning domain can be seen in figure 2 .

There are three types of machine learning
To improve validation accuracy with less number of epochs Using CNN model.

- Supervised Machine Learning
- Unsupervised Machine Learning
- Reinforcement Machine Learning

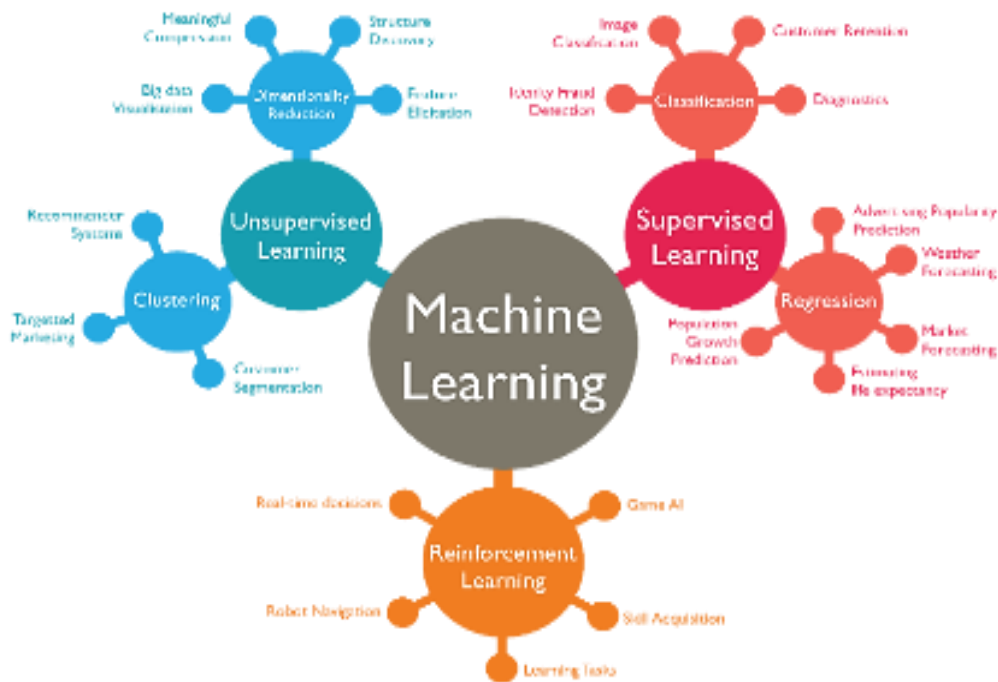


Figure 2: Structure of Machine Learning domain

II. LITERATURE REVIEW

Angeline Christopher. Y and Dr. Sivaprakasam [4], utilizing decision tree classifiers (CART) on cancer databases, obtain sensitivity of 69.23%. A. Pradesh compares the efficacy of the machine learning strategies SVM, NN, and BF Tree. SMO's accuracy is superior to that of other predictors, according to the data. Using Wisconsin

cancer (original) databases and neuron fuzzy approaches, Joe achieves an accuracy of 95.06 percent. In this study, a hybrid approach was suggested to improve the Wisconsin cancer (original) sets' classifier (95.96%) with a ten-fold cross-validation approach.[10]

In order to predict basal cell endurance, Liu Ya-Qin, W. Cheng, and Z. Lu[11] performed experiments on osteosarcoma data using the C5 classifier with classifiers.

They did this by providing substantial dataset to train based on the initial set using permutations with sequences to source multiple sets of data that are the same size as the first records. The records of 202,932 terminally ill patients are taken by Delen89 et al. Lu19 and pre-classified into two categories: "survived" (93,273) and "not lived to tell the tale" (109,659).[3]

III. PROBLEM FORMULATION

- Real-world 3D objects should be recognisable from 2D photos using computer vision.
- There are available values for r, g, and b. We now need a new technique to separate the color name from the RGB value.
- We calculate a distance (d), which shows how close we are to selecting the choice with the least distance, before choosing the color name.

IV. OBJECTIVES

- To Study a color from a photograph taken with a camera and to test the predictions generated by the machine learning system under various lighting conditions.
- To enhance the performance of the model by increasing the epochs.

V. METHODOLOGY

This work makes use of the Wisconsin Breast Cancer (original) datasets from the UCI Machine Learning Repository. 569 cases (350 innocuous and 219 malignant), 2 classes (with 65.5 percent malignant and 34.5 percent benign), and 11 integer-valued characteristics are included in the data set.[12]. Figure 3 shows the flow chart of the model with open CV

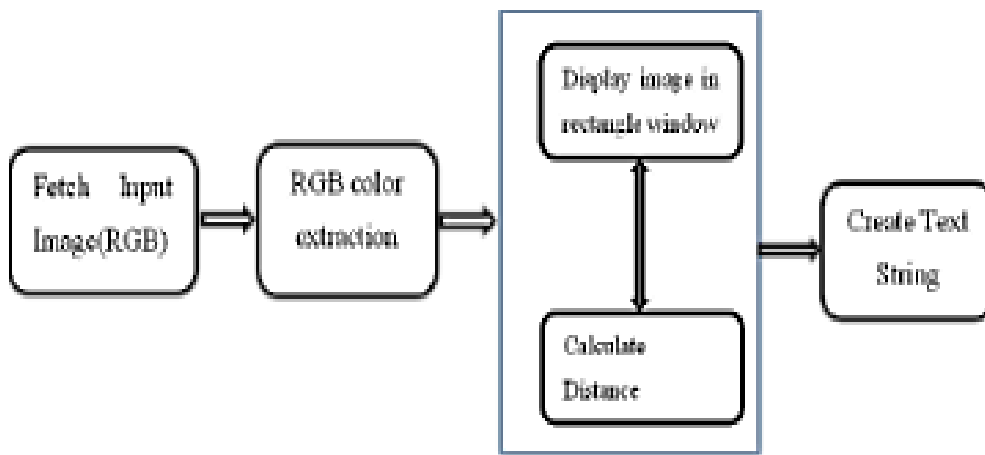


Figure 3: Flow chart of the model with open CV

A. Dataset Attribute Information

There are different attributes in dataset of "Diagnostic Wisconsin breast Cancer Database". Such attributes include:

- 1) Id
- 2) "Diagnosis" (B = benign; M = malignant)
- 3) For each centrosome, ten characteristics will be calculated, including:

- Radar
 - The tone
- Fence line
- Territory
- Easiness
- Permeability (area - 1.0 / perimeter 2)
- Curvature
- Concave edges
- Symmetry
- Dimensional fractals ("coastline approximation" - 1)
- 4) The three sections of attributes 3-32 are as follows:
- Mean (3-13),

- Stuck-Error (13-23)
- Worst (23-32)
- 5) Each has ten parameters. ("Radius", "texture", "area", "perimeter", "smoothness", "compactness", "concavity", "concave points", "symmetry" and "fractal dimension")
- Mean
- Common Phrase
- Stupidest

In order to diagnosing breast cancer, we must make sure our model does not create too many false-positives. For instance, (Anyone don't have cancer, but we told them to go for the treatment) or false-negatives (someone has cancer, but we have told hat individual not to for the treatment). That is the main cause, the highest overall accuracy model is chosen.

B. Co-Relation Map between Attributes

Correlation mapping on the code snippet can be seen in figure 4. Feature relation map can be seen in figure 5

```

In [8]: #correlation map
f,ax = plt.subplots(figsize=(18, 18))
sns.heatmap(Var_x.corr(), annot=True, linewidths=.5, fmt= '.1f',ax=ax)
  
```

Figure 4: Correlation map on the model code

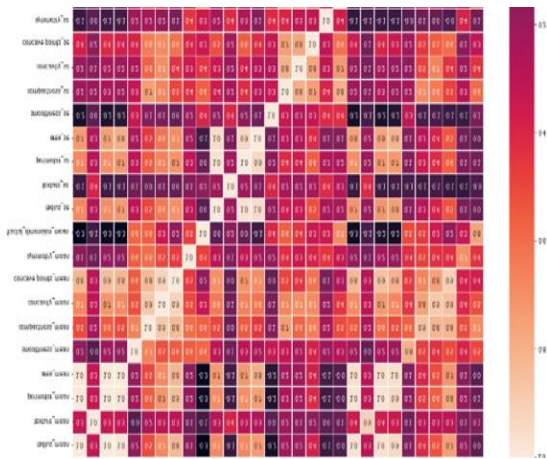


Figure 5: Feature relation map

C. Co-Relation Map Observation

If we look the Map, we can immediately find out the co-relation between some parameters. For instance, the radius-mean got co-relation having the perimeter-mean and area-mean values, respective, of 1 and 0.99. Because the physical size of the observation may theoretically be the same info for all three from these characteristics, I must choose one of the six because of these parameters in order to go into further analysis.

After that take part of a cell as roughly from a circle, then the formula for its radius should be its r . Its perimeter and area are then $2\pi r$ and πr^2 , respectively. A cell's radius is the crucial of its size. That is why, it is essential to select radius as our attribute to represent the size of a cell.

Similarly, there is co-relation exist between the parameter's compactness, concavity, and concave points same as the size parameters. Take one of these three parameters that got the data about the shape of the cell.

Compactness is an parameter name I am going to discard the other two parameters. [13]

According to the dataset description, Parameters like radius, perimeter area are perfectly understandable to me. However, what do texture, What do homogeneity, flatness, flatness, curvature, and geometric scales actually mean?

The test statistic of gray-scale is referred to as smoothness. Each pixel of an image is represented by the 8-bit integer, or a byte, from 0 to 255 providing the amount of light, where 0 is clear black and 255 is clear white. The darker the image is the lower is the mean of intensity level of a pixel, i.e. byte. So, the SD of gray-scale values means how intense levels are spread for particular individual cells. The higher SD the more contrasting the image is. [14]

Next, smoothness is measured as the disparity between it radially line's strength and the normal length of the two real axis around it.

The contouring is straight in that area if the number is low.
$$\text{Smoothness} = \sum_i ||r_i - r_{i-1} + r_i + 12|| \text{perimeter}$$

The concavity is captured by drawing chords between two boundary points, which lie outside the nuclear. For the concavity-mean the mean value of these lengths is calculated.

To calculate the main axis through the center, asymmetry, is discovered. Then, assess the distance differences in both directions with lines parallel to the nuclear border and the main axis.

Using successively bigger rulers, the circle of the center is examined in order to estimate the spatial structure; as the ruler size grows, the perimeter shrinks. We may determine the entropy by measuring the downslide and plotting the regression analysis of the ruler size to against log of the perimeter. A greater number for any of the form characteristics indicates a less regular contour and, thus, a higher likelihood of malignant. [15]

D. Feature Selection

Feature selection mapping on the code is shown in figure 6 and feature properties can be seen in figure 7 .

```
In [9]:
#Feature Selection

# first ten features

data_dia = Var_y
data = Var_x
data_n_2 = (data - data.mean()) / (data.std()) # standardization
data = pd.concat([Var_y, data_n_2.iloc[:,0:10]], axis=1)
data = pd.melt(data, id_vars="diagnosis",
               var_name="features",
               value_name="value")

plt.figure(figsize=(10,10))
sns.violinplot(x="features", y="value", hue="diagnosis", data=data, split=True,
               inner="quart")
plt.xticks(rotation=90)
```

Figure 6: Feature selection mapping on the system code

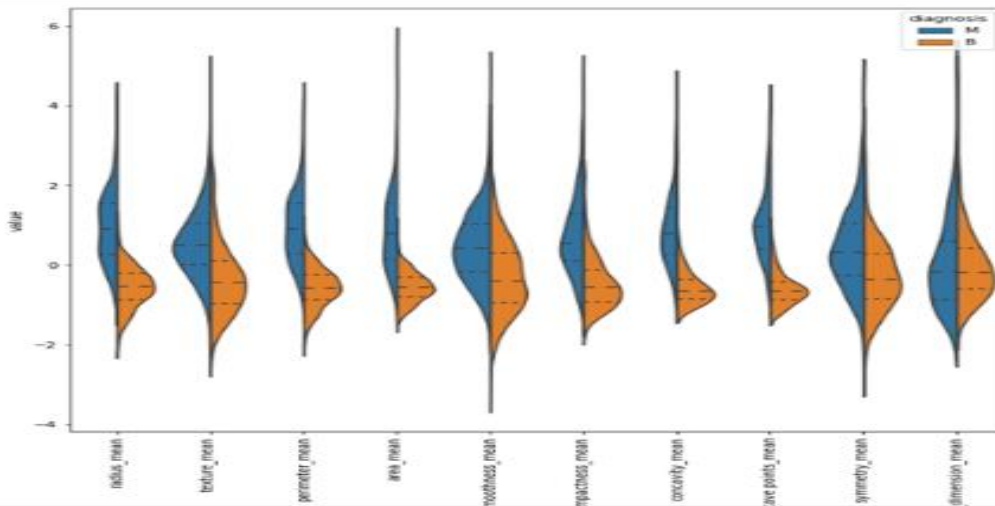


Figure 7: Feature properties

- *Co-Relation Map of First 10 Features*

The correlation features can be seen in figure 8

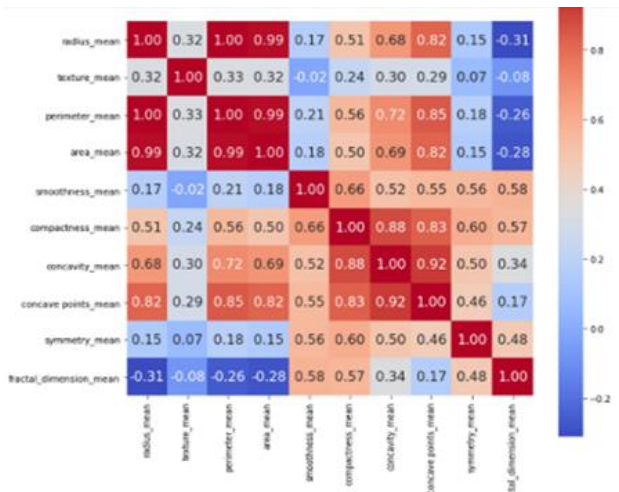


Figure 8: Correlation between features

I have found a few fascinating linear designs. The nearly flawless geometric patterns, for instance, between its radius, periphery, and area values show that there is a correlation between all these factors. added variables could have

possibly co-relation are the concavity, concave-points and compactness.

- *Co-Relation Map of First 10-20 Features*

Figure 9 shows the feature correlation.

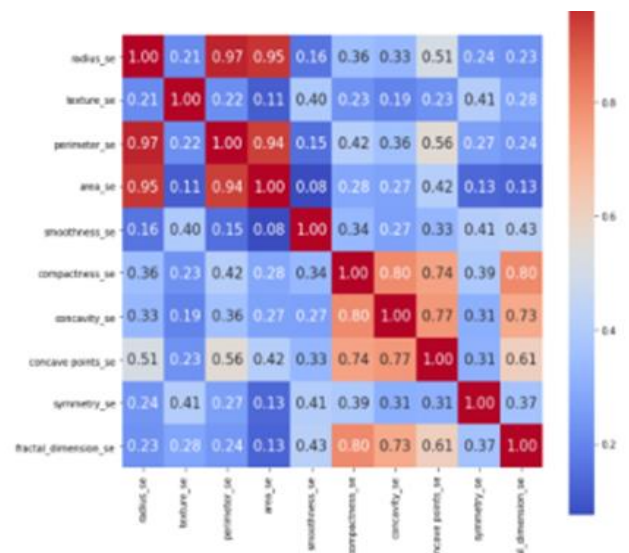


Figure 9: Feature correlation

- *Second Ten Features*

The feature properties can be shown in figure 10

```
In [14]:
# Second ten features
data = pd.concat([Var_y, data_n_2.iloc[:,20:31]],axis=1)
data = pd.melt(data,id_vars="diagnosis",
               var_name="features",
               value_name="value")

plt.figure(figsize=(10,10))
sns.violinplot(x="features", y="value", hue="diagnosis", data=data,split=True, inner="quart")
plt.xticks(rotation=90)

Out[14]:
```

Figure 10: Feature properties

- *Co-Relation Map of First 10-20 Features*

Figure 11 shows the correlation between the properties

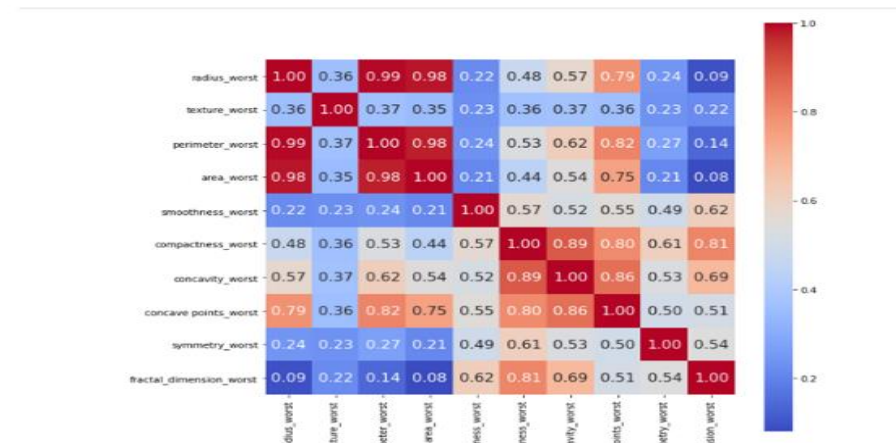


Figure 11: Correlation between properties

- *Patients' Counts M or B type Cancer*

Figure 12 shows the Cancer type distribution

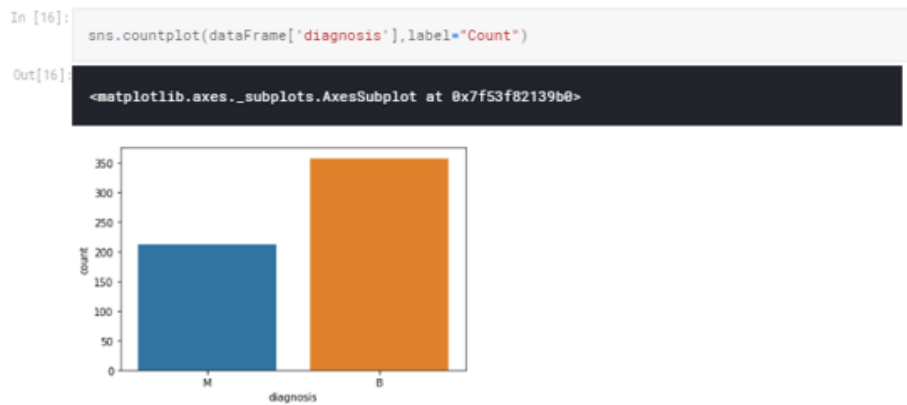


Figure 12: Cancer type distribution

- *Linear Patterns Between Highly Co-Related Features*

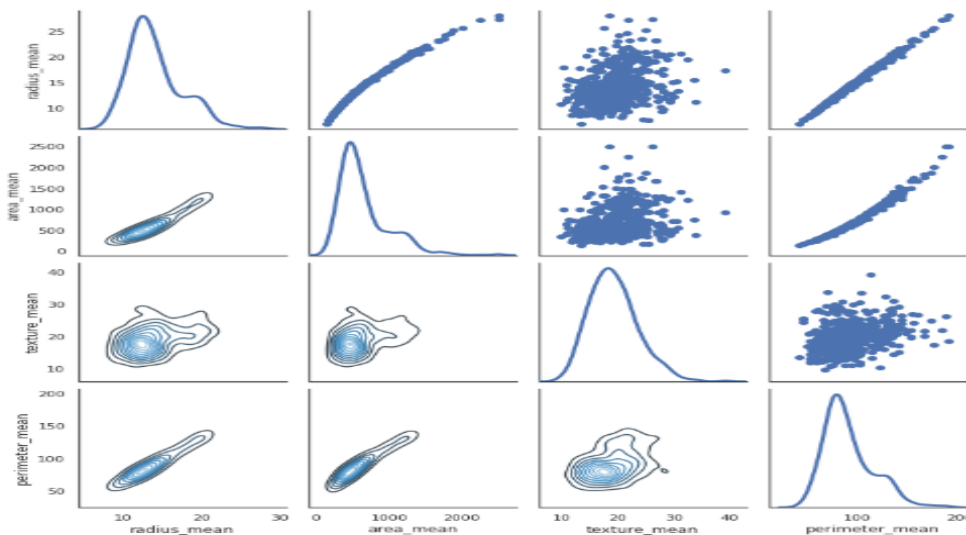


Figure 13: Distribution patterns

I have found a few fascinating linear designs. Distribution model can be seen in figure 13. The pretty perfect exponential sequences, for instances, in between radius, periphery, and area values show that there is a correlation between these factors. added values could have possibly

co-relation are the concavity, concave-points and compactness. Figure 14 and 15 show the variety of parameters used in training the models

1) Selected Parameters for Training the Model

	texture_mean	area_mean	smoothness_mean	concavity_mean	symmetry_mean	fractal_dimension_mean	texture_se	area_se
0	10.38	1001.0	0.11840	0.3001	0.2419	0.07871	0.9053	15
1	17.77	1326.0	0.08474	0.0869	0.1812	0.05667	0.7339	74
2	21.25	1203.0	0.10960	0.1974	0.2069	0.05999	0.7869	94
3	20.38	386.1	0.14250	0.2414	0.2597	0.09744	1.1560	27
4	14.34	1297.0	0.10030	0.1980	0.1809	0.05883	0.7813	94

Figure 12: Parameters from texture mean to texture size

se	area_se	smoothness_se	concavity_se	symmetry_se	fractal_dimension_se	smoothness_worst	concavity_worst	symmetry_worst
	153.40	0.006399	0.05373	0.03003	0.006193	0.1622	0.7119	0.4601
	74.08	0.005225	0.01860	0.01389	0.003532	0.1238	0.2416	0.2750
	94.03	0.006150	0.03832	0.02250	0.004571	0.1444	0.4504	0.3613
	27.23	0.009110	0.05661	0.05963	0.009208	0.2098	0.6869	0.6638
	94.44	0.011490	0.05688	0.01756	0.005115	0.1374	0.4000	0.2364

Figure 13: Parameters from smoothness to symmetry

VI. SYSTEM ARCHITECTURE

A. Dataset

In this study, the Wisconsin Breast Cancer dataset from the Machine Learning Repository was used. The information has 30 numerical characteristics descriptors. From a digital picture of a lumbar puncture aspirate (FNA) of a breast

mass, features are estimated. They characterize the attributes of the visible cell nuclei in the picture. 569 cases total in the dataset are categorized into benign and tumorous classifications. There are 212 benign cases and 357 invasive cases in all (UCI, 2019). [16]

B. Experiment Method

Figure 16 displays the test's architecture.

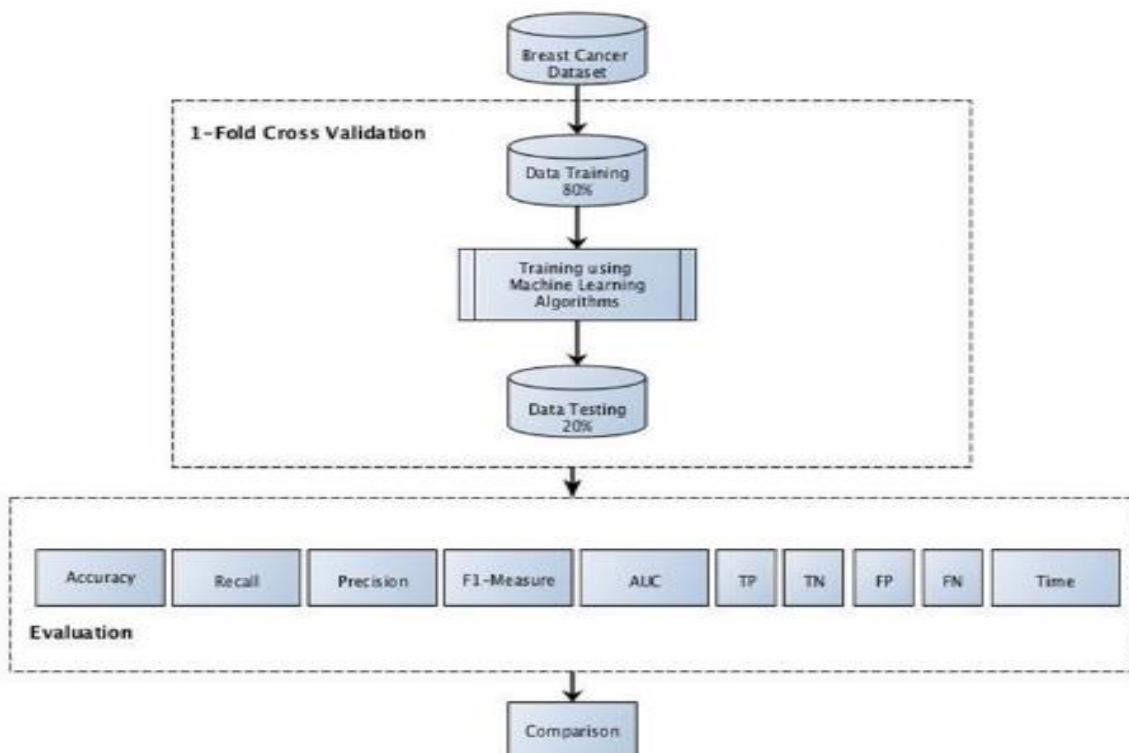


Figure 16: Cross Confirmation of Folds

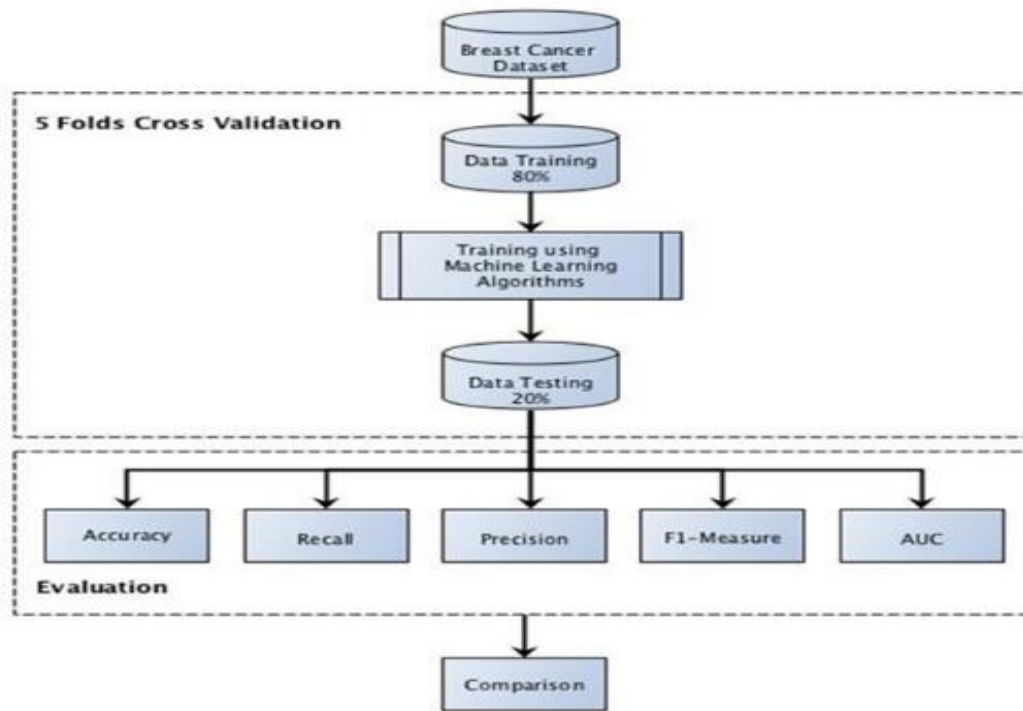


Figure 17: Working in Five Folds

Figures 16 and 17 demonstrate how we tested the effectiveness of data mining algorithms in two different ways: (1) One-fold discriminant analysis; and (2) five-fold cross validation. We employed the performance measures F1-Score, Accuracy, Recall, Precision, and Area Under Curve (AUC). We divided the data in this trial into data sets (80%) and testing data (20%). Next, we run each algorithm through one round of testing and training (1-fold cross validation). After that, we do five rounds (5-fold merge) utilizing various, randomly generated data from the data set for both the testing and the training phases. We then average each performance index. On a desktop computer macOS High Sierra with a 2.5GHz Intel Core i5 and 16 GB RAM, the entire research is carried out using Python 3.7..

C. Confusion Matrix and Quality Measures

First, true positive, true negative, false positive, and false negative definitions are provided. [17]

The following is a description of the benchmarks that we employed.

- *Accuracy*

Accuracy is a measure of how accurately the model was trained. It is described as the algorithm for measuring quality by comparing right judgments to all other guesses.

- *Recall*

The ratio of accurately identified positive cases to True Positives and False Negatives is known as recall. The recall function and recall Table 1 is shown below

$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

- *Precision*

In calculating the ratio among both True Positives and all positive predictions, precision refers to the degree of

accuracy. Below is a presentation of the fineness calculation:

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

Table 1: Recall Table

	Class = True	Class = False
Class = True	True positive	False Negative
Class = False	False Positive	True Negative

- *F1-Score/F-Measure*

A weighted average of Precision and Recall is the F1-Score. Listed below is a presentation of the F1-Scores eqs:

$$\text{F1-Score} = \frac{2 * (\text{Precision} * \text{Recall})}{(\text{Precision} + \text{Recall})}$$

- *AUC*

By considering the number of the area beneath the ROC curve, AUC calculates competence. An algorithm has highest efficiency if its area under the ROC curve is larger.

D. K-Fold Cross Validation

A statistical method for verifying and assessing learning algorithms or models is called cross-validation. Through the process of crossvalidation, data is randomly divided into training and testing sets. The simplest type is k-fold cross validator, which uses one of the k divisions as a validation set. To ensure that our measures of success accurately reflected the whole dataset, we employed k-fold merge (Suyanto. 2018). [18]

- Machine Learning Algorithms

We will go over every engine we utilized in this investigation in this part. Each fundamental algorithmic idea with regard to the classification issue:

- Evaluation of the Algorithms Using One-Fold Stepwise Regression

We initially used the 1-fold test set approach to compare the systems.

- Analysis of the Solutions Using a Fivefold Validation Method

We used the five-fold cross validation approach, and we used the average single value. Results for the second-best method, KNN, are reported in Table 2. The computational time of XG Boost, which is faster than other programs while being cheaper than GNB, is 0.08 seconds. GNB showed improvement and processes in 0.008 seconds. Surprisingly, the SVM approach, one of the top algorithms in earlier research, takes the longest to compute and yields one of the worst performance measure scores.[20] XG Boost provides the best outcome when compared to other models, which is similar to the prior finding. The performance metrics for ANN include recall 96,75 percent, precision 97,28 percent, F1-Score 96,99 percent, and accuracy 97,19 percent. XG Boost also has the highest AUC at 99,61 percent. The studies suggest that ANN is the technology that runs the strategies for identifying breast cancer in the Wisconsin Breast Cancer dataset. [19]

VII. SIMULATION AND RESULTS

From all the comparative analysis of 7 models which include SVM, AdaBoost, Random forest, KNN, XG Boost, Bagging and ANN the data training dataset count was slowly increased from 0.1 to 1.0 and the model parameters which include training Accuracy, Testing Accuracy, Validation Accuracy, F1 score, Recall, and Precision were recorded and saved

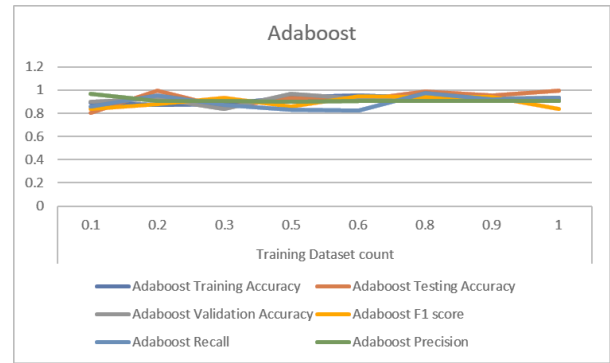


Figure 20: Adaboost Model Analysis

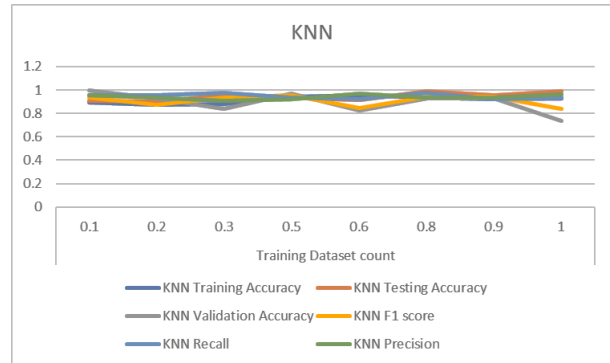


Figure 21: KNN model analysis

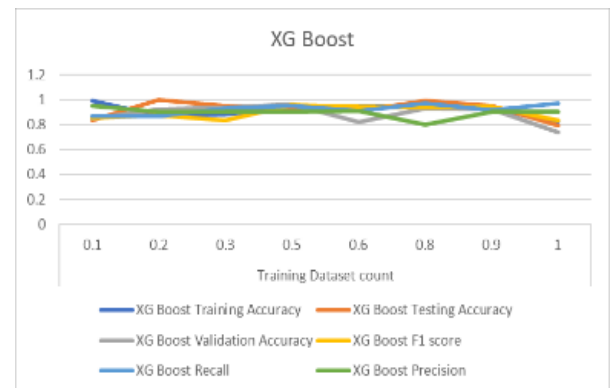


Figure 22: XG Boost Model Analysis

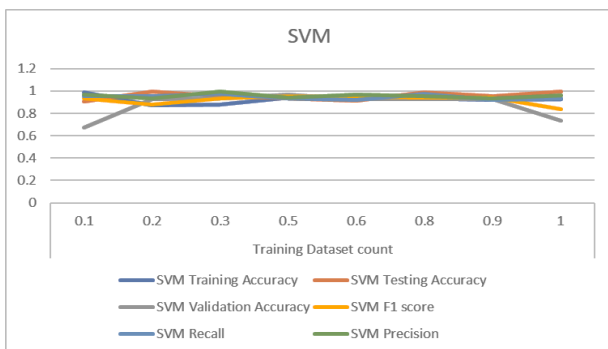


Figure 18: SVM Model Analysis

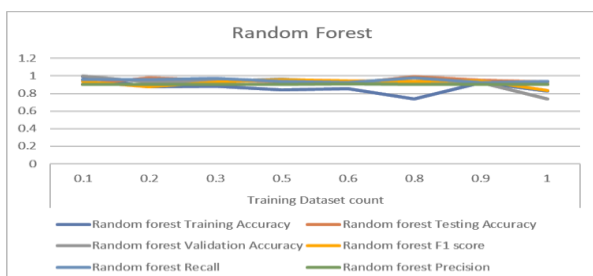


Figure 19: Random Forest Model Analysis

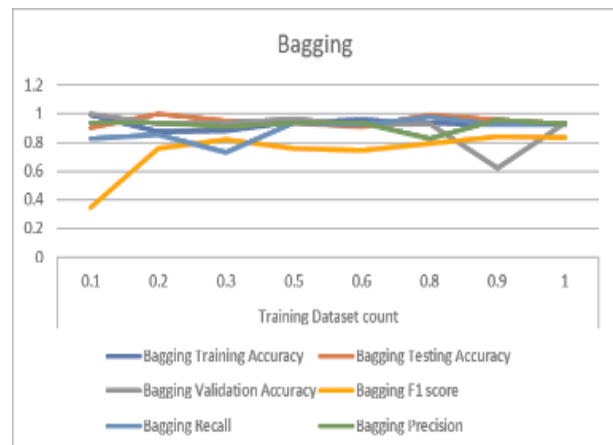


Figure 23: Bagging Model Analysis

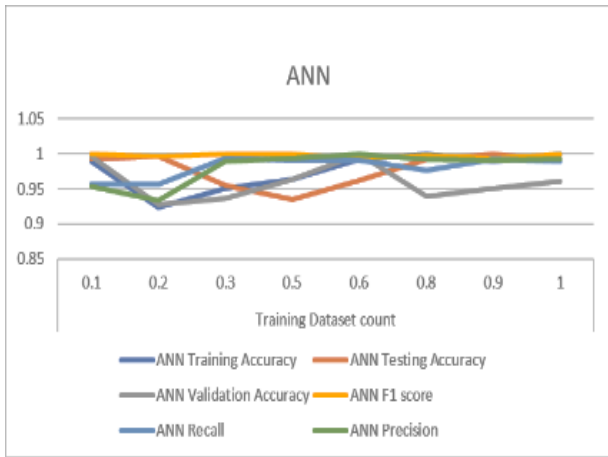


Figure 24: ANN Model Analysis

A. Comparative Analysis

Further, it was observed that the SVM model’s finest parameters were recorded as Highest training accuracy = 0.9916, the highest testing accuracy = recorded to be 0.97, the highest validation accuracy 0.92, F1 score was 0.94, the Recall was 0.97 and the least Precision was recorded as 0.925 as shown in figure 18

For random forest it was observed that Training accuracy = 0.99, Testing accuracy = 0.97, Validation accuracy = 0.96, F1 score = 0.93, Recall = 0.97, Precision= 0.903 as can be seen in figure 19

For AdaBoost the highest and the finest parameters were recorded as Training accuracy= 0.95, Testing accuracy= 0.99, Validation accuracy=0.96, f1 score = 0.945, Recall = 0.93, Precision= 0.9026 as seen in figure 20.

For KNN it was observed that Training Accuracy= 0.95, Testing accuracy = 0.95 Validation Accuracy=0.9 F1 score=0.97 Recall=0.97, Precision=0.905 as seen in figure 21

For XG boost, Training Accuracy= 0.99, Testing Accuracy= 0.93, Validation Accuracy= 0.96, F1 score=0.95, Recall= 0.97, Precision=0.94 as shown in figure 22.

For Bagging it was observed Training Accuracy=0.99, Testing Accuracy=0.95, Validation Accuracy=0.99, F1 score=0.85, Recall=0.97, Precision=0.99 as shown in figure 23

After observing and comparing all the models under all the different parameters, it was observed that ANN outperforms all the models in all the parameters in the comparative analysis. With the highest and the finest parameters reported as Training Accuracy= 1.0, Testing Accuracy=1.0, Validation Accuracy=1.0, F1 score=0.999, Recall=0.993, Precision=1.0.

This result comparison is compiled in Table 2. Figure 25 shows the comparative analysis of the model

Table 1: Comparison table

	SVM	Random Forest	Ada Boost	KNN	XG Boost	Bagging	ANN
Training Accuracy	0.99	0.97	0.99	0.95	0.99	0.99	1
Testing Accuracy	0.97	0.96	0.96	0.95	0.93	0.95	1
Validation Accuracy	0.92	0.96	0.94	0.9	0.96	0.99	0.98
F1 score	0.94	0.93	0.93	0.97	0.95	0.85	0.99
Recall	0.97	0.97	0.93	0.97	0.97	0.97	0.97
Precision	0.92	0.9	0.9	0.9	0.94	0.93	1

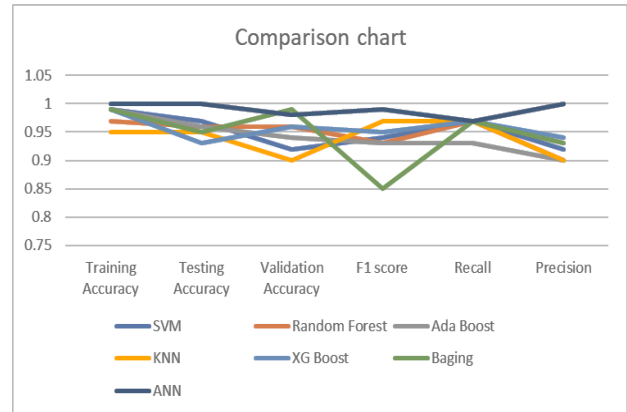


Figure 25: Comparative analysis of the models

VIII. CONCLUSION

In this study, we used Employing eight distinct algorithms, the Wisconsin Breast Cancer dataset is used to categorize pancreatic cancer. In order to determine the performance index, we used 1-fold and 5-fold merge. Optimum method. The analysis revealed that ANN has the greatest AUC of 99.99 percent and the best performance metric when compared to other algorithms. We conclude that utilizing the Wisconsin Breast Cancer dataset, ANN is the best accurate algorithm for classifying ovarian cancer. Future evaluations of ANN versus computation that weren’t used in this experiment and on various datasets. The results analysis shows that the combination of multidimensional data with various feature selection, classification, and dimensionality reduction approaches might offer advantageous tools for inference in this field. It is necessary to do further study in this area to improve the classification systems' performance so that they can make predictions about more factors.

CONFLICTS OF INTEREST

The authors declare that they have no conflicts of interest.

REFERENCES

- [1] Key, T. J., Verkasalo, P. K., & Banks, E. (2001). Epidemiology of breast cancer. The lancet oncology, 2(3), 133-140.
- [2] 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 and Prevention, and National Cancer Institute; 2012.
- [3] Chaurasia, V., & Pal, S. (2014). Data mining techniques: to predict and resolve breast cancer survivability. International

Journal of Computer Science and Mobile Computing
IJCSMC, 3(1), 10-22.

- [4] Djebbari, A., Liu, Z., Phan, S., & Famili, F. (2008). An ensemble machine learning approach to predict survival in breast cancer. *International journal of computational biology and drug design*, 1(3), 275-294.
- [5] Aruna, S., Rajagopalan, S. P., & Nandakishore, L. V. (2011). Knowledge based analysis of various statistical tools in detecting breast cancer. *Computer Science & Information Technology*, 2(2011), 37-45.
- [6] Agarap, A. F. M. (2018, February). On breast cancer detection: an application of machine learning algorithms on the wisconsin diagnostic dataset. In *Proceedings of the 2nd international conference on machine learning and soft computing* (pp. 5-9). V. Chaurasia, S. Pal, and B. Tiwari, "Prediction of benign and malignant breast cancer using data mining techniques," *Journal of Algorithms & Computational Technology*, vol. 12, no. 2, pp. 119-126, 2018.
- [7] Fatima, N., Liu, L., Hong, S., & Ahmed, H. (2020). Prediction of breast cancer, comparative review of machine learning techniques, and their analysis. *IEEE Access*, 8, 150360-150376.
- [8] Toprak, A. (2018). Extreme learning machine (elm)-based classification of benign and malignant cells in breast cancer. *Medical science monitor: international medical journal of experimental and clinical research*, 24, 6537.
- [9] Jacob, D. S., Viswan, R., Manju, V., PadmaSuresh, L., & Raj, S. (2018, March). A survey on breast cancer prediction using data mining techniques. In *2018 Conference on Emerging Devices and Smart Systems (ICEDSS)* (pp. 256-258). IEEE.
- [10] Padhi, T., & Kumar, P. (2019, January). Breast Cancer Analysis Using WEKA. In *2019 9th International Conference on Cloud Computing, Data Science & Engineering (Confluence)* (pp. 229-232). IEEE.
- [11] Ren, Johnny. Investigation of Convolutional Neural Network Architectures for Image based Feature Learning and Classification (2016)
- [12] Bobulski, J., and Kubanek, M - Waste Classification System Using Image Processing and Convolutional Neural Networks. Springer International. (2019)
- [13] Mandar Satvilkar – Image-Based Trash Classification using Machine Learning Algorithms for Recyclability Status (2018)
- [14] Gaurav Mittal, Kaushal B. Yagnik, Mohit Garg & Narayanan C. Krishnan – SpotGarbage: a smartphone app to detect garbage using deep learning. (2016).
- [15] Rad, Kaenel, Droux - A computer vision system to localize and classify wastes on the streets (2017).
- [16] Piotr. N and Teresa. P - Application of deep learning Color classifier to improve e-waste collection planning (2020) 1-9.
- [17] Ashok.S, Arun Kumar.HN, Niranjana, Shekar.A, Arun Kumar.DR - A Deep Learning Approach for Real-Time Garbage's Detection and Cleanliness Assessment (2020).
- [18] P. Sermanet, D. Eigen, X. Zhang, et al. .Overfeat: Integrated Recognition, Localization and Detection Using Convolutional Networks. ICLR, 2014
- [19] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet Classification with Deep Convolutional Neural Networks. NIPS, 2012