

# Classification of Breast Cancer Data Using Enhanced Supervised Machine Learning Algorithm Optimized SVM

N. Aksha<sup>1</sup>, Dr. Naresh Alapati<sup>2</sup>

<sup>1</sup> M.Tech Student, <sup>2</sup> Assoc. Professor

Department of Computer Science & Engineering

Vignan's Nirula Institute of Technology and Science for Women, Palakaluru

[chinupapai803@gmail.com](mailto:chinupapai803@gmail.com)<sup>1</sup>, [alapatinaresh13@gmail.com](mailto:alapatinaresh13@gmail.com)<sup>2</sup>

## Abstract

Breast cancer has been a hot topic in health care informatics for the last few years, as it is the second major cause of cancer-related deaths in women. Breast cancer can be identified by a biopsy where tissue is removed and studied under a microscope. The diagnosis is based on the qualification of the histopathologist to look for abnormal cells. However, if the histopathologic is not well trained, this may lead to a misdiagnosis. With recent advances in the analysis of medical data and machine learning, there is an interest in trying to develop reliable pattern recognition based systems to improve the quality of diagnosis. This paper focuses on the classification of breast cancer as binary labeled data because it is benign or malignant. The objective is therefore to determine whether breast cancer is benign or malignant and to predict the recurrence and non-recurrence of malignant cases after a certain period of time. To achieve this, we used machine learning classification methods to fit a function that could predict a discreet class of new inputs. For this purpose, the enhanced SVM-Support Vector Machine (SVM) is a supervised machine learning algorithm that can be used for both classification and regression analysis. It separates two classes using a hyperplane.

**Keywords:** Breast Cancer, Classification, Machine Learning Algorithm, Supervised Learning, SVM Classifier, Feature selection, Disease Prediction.

## 1. Introduction

Breast cancer (BC) is the most widely recognized cancer in women, affecting around 10% of all women in certain phases of their lives. At present, the rate continues to increase and information shows that the

endurance rate is 88 per cent after five years of determination and 80 per cent after 10 years of analysis. Early predictions of breast cancer have so far led to progress, increasing the pace of breast cancer by 39%, starting in 1989. Due to the shifting nature of breast cancer side effects, patients are often exposed to a torrent of tests, including but not limited to mammography, ultrasound and biopsy, to check their likelihood of developing breast cancer. Biopsy is the most characteristic of these methods, including the extraction of test cells or tissues for assessment. The example of the cells is obtained from the breast fine needle longing (FNA) methodology and then sent to the pathology lab to be analyzed under a magnifying lens[7].

Numerical highlights, such as sweep, surface, edge and zone, can be estimated from minute pictures. Information later acquired from FNA is broken down in combination with different imagery information to predict the likelihood of a patient having a harmful breast cancer tumor. In this situation, the computerized framework would be extremely valuable. It is likely to facilitate the procedure and improve the accuracy of the specialist's forecasts. Furthermore, whenever the wealth dataset and robotic framework are properly maintained, it will conceivably dispense with the requirements for patients to undergo different tests, such as mammography, ultrasound, and MRI, which subject patients to an enormous amount of agony and radiation. Overall, the early forecast remains one of the indispensable points of view in the subsequent procedure.

Information mining strategies or characterization can reduce the number of false positive and false negative choices. As a result, new ways such as disclosure of

information in databases (KDDs) have become a preferred device for restorative analysts. In this paper, six order models were used; Decision Tree, K-Neighbors, Linear Discriminant Analysis (LDA), Logistic Regression, Naïve Bayes and Support Vector Machine (SVM) were run on Wisconsin Breast Cancer (Unique) datasets, both when applying Principal Component Analysis. The results obtained are then estimated using different execution measurements to contrast the calculations with the discovery of the most appropriate cancer expectation model.

Machine learning is a section of artificial intelligence that has a place in the science and design of making smart machines. Mechanized information acquisition by machine learning through the structure and use of calculations, where accurate information is required by calculations. Essentially, machine learning procedures are trained by machine learning that relies on the use of probability. There are distinctive ways to have a place in machine learning. Administered learning: In directed learning, starting with datasets containing the preparation of models that can be recognized by the associated level. It does this by running information through the calculation of learning. The objective of administered learning is to effectively recognize the new information provided through managed learning and to use past information collection and learning calculations to gain knowledge of the information-distinguishing strategy. The calculations under managed learning take the information sources that the yield is now known for explanation all together that the calculations will make the machine to discover by holding it in contrast to the specific yield and the definitely realized yield to test for any further blunders. The machine is therefore shown.

Well-known administered learning calculations include classification, increase in inclination,

expectation and relapse. The model is therefore adjusted at that point. With such calculations, the machine makes use of managed learning to attempt, by misuse of proper examples, to make the prediction of mark esteems on unmarked data. Managed learning finds the machine at any location in such zones, the more out-of-the-box opportunities are normal through recorded data. Solo learning: Unsupervised learning ponders anyway frameworks will figure out how to speak to explicit info designs in a way that reflects the applied math structure of the collection of information designs. By distinguishing from directed learning or fortification learning, there are no explicit target yields or natural assessments identified with each info; rather, the unattended student is in contact with past inclinations as to what parts of the information structure should be caught within the yield. A particular output is not to be made by solo learning.

The learning specialist points to the discovery of the structures and examples in the information. Semi-supervised learning: under this type of machine learning, the machine is designed to be suitable for learning every named and untagged data for instructive reasons. This includes, in particular, the preparation of the machine by means of a small amount of marked information to the detriment of a larger than usual amount of untagged information. This may be the rationale method by which untagged data are practical and clearly to be collected. In many cases, this type of machine learning is used for calculations such as classification, forecasting and relapse. In addition, this kind of learning is used within the field where the cost of a related naming is spent too much to do because of a fully-named training strategy. Recommended use of semi-directed learning is a face recognition through a computerized camera.

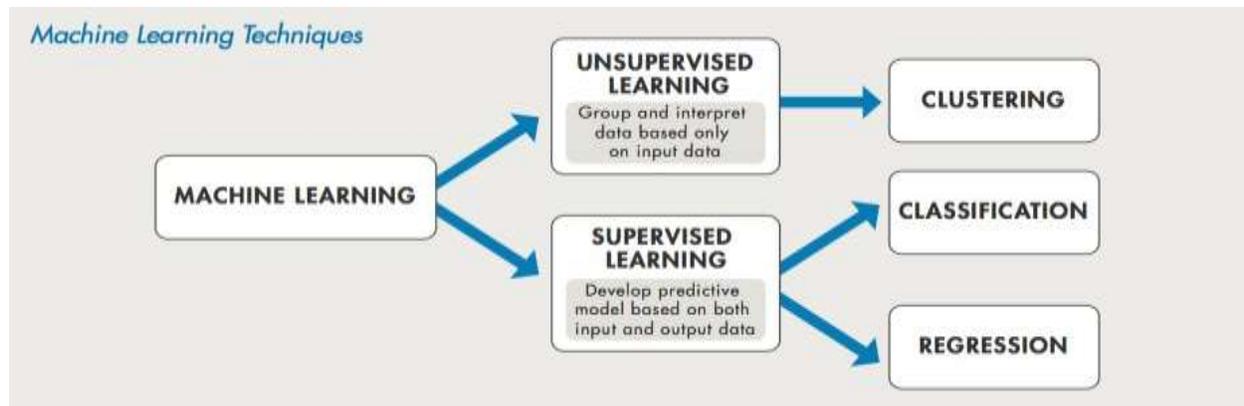


Fig. 1 Machine Learning Techniques

Supporting learning: under this type of machine learning, machine learning calculations go through the experimentation method to deal with the positive structure of the activities that deliver the most straightforward outcomes, and to develop applications within the field of play, route, and man-made reasoning. It is generally used for computerized reasoning, gaming, and route. There are three components that operate mainly under this type of machine learning—the specialist, the student, the air with which the operator communicates and, furthermore, the activities that the specialist intends to undertake. The aim of supporting learning is to shape a specialized selection of activities that will encourage increased reward over an ideal measure of time. In this way, the arrangement is clear that the fortification enables the machine to familiarize itself with the easiest strategy to come up with in order to achieve the best results. Community Learning: Recommendations are developed by means of a method known as synergistic separation, which is an essential type of recommended framework. Among the vast number of decisions, and depending on the examination of client inclinations, it encourages clients to discover something of importance.

## 2. Literature Survey

Agglomeration or clustering means are intended to extract data from the information index in order to obtain groups or clusters and to represent the data set. Classification otherwise referred to as supervised learning in machine learning, intends to characterize

obscure things by supporting the learning of existing examples and classes in the data set and by anticipating future things. The preparation set, which is used to assemble the grouping structure, and therefore the investigation set, which will generally evaluate the classifier, is commonly referred to in the classification assignments[2]. In addition, basic advancements in breast cancer survival expectations have been completed using marked, unmarked, and pseudo-named quiet information. Prognosis of breast cancer survival has been helped by machine learning calculations that can predict the endurance of a specific patient who is dependent on chronic tolerant information. Neural systems and related strategies are very committed to the anticipation of breast cancer. In recent decades, the Artificial Neural Networks have been increasingly used by an increasing number of specialists and have become a functioning examination zone. The ANNs have achieved various victories with incredible advances in Breast Cancer classification and analysis in the early stages.

Furundzic et al.[2] attempted to establish a neural framework for the assurance of the threat to the chest. Negative association was used to get ready for computing, therefore, to rot the issue and get it. In this article, the author discussed two procedures, e.g. the formative strategy and the outfit approach, in which transformative philosophy can therefore be used to lay down the neural framework. The outfit approach was expected to address the huge issues that were still in progress.

Floyd et al [3] carried out an audit of data digging frameworks for quality decision-making. This article

oversaw the most widely used data burrowing techniques for quality decision-making and threatening development requests, focusing in particular on four basic rising fields. They are neural framework-based computations, AI estimates, inherited figuring and pack-based computations and have decided on future improvements in this field.

Fogelet al .[4] developed a mechanized chest disease end by joining the inherited count and back spread neural framework, which was designed as a faster classifying model to reduce break down time in the same way as extending the precision of mass collection in the chest to either compassionate or compromised. In these two different clearing structures, the data set was finished. In Set An, only records with missing characteristics were cleared, while Set B was set up with a common true-to - life cleaning technique to recognize any ominous or missing characteristics. At last Set A gave 100 percent of the highest precision rate, and Set B gave 83.36 percent accuracy. Thus, the manufacturer has assumed that remedial data is best maintained in its exceptional motivating force as it provides a high precision rate when it differs from balanced data.

Fogelet al .[5] analyzed the creation of neural frameworks for the perception of chest threat and related works for the assurance of chest disease using a multilayer perceptron backspreading technique. Rather than back spread, The manufacturer used 699 data, which had missing characteristics and was removed , leaving 683 data. Two test plans have been conducted using these characteristics. The main breakup consisted of five fundamentals with 9-2-1 Multi Layer Perceptron (i.e. 9 data, 2 masked center

points and 1 yield center point) and a second investigation involving 9-9-1 Multi Layer Perceptron. The result of the main test after 400 years in each of the five fundamentals was 97.5 per cent accurate. In the second study, in relation to the previous assessment, the best execution was discovered at a precision rate of 98.2 percent for the lesser hidden center points.

### 3. Proposed Model

The general research methodology for this assessment was balanced against the background of the data disclosure process. The data verification stage was the main time of this way of thinking, in which we obtained the relevant data for the examination. The next stage was the pre-processing phase of the data in which the accumulated information was compiled, cleaned and modified, with the ultimate objective that the datasets were sensitive to the request. After that, still in the next stage, we finished the extraction of the part.

#### Support Vector Machine

Support Vector Machine ( SVM) is a regulated algorithmic machine learning standard that can be used for any classification or relapse difficulty. However it may be, it is essentially used in classification issues. In this algorithmic standard, every data item is plotted as a point in ndimensional space where n is the number of highlights one has, with the estimation of each element being the estimation of a specific arrangement [3]. At that point, we perform classification by finding the hyperplane that separates the two classes all around, which is shown in the figure below:

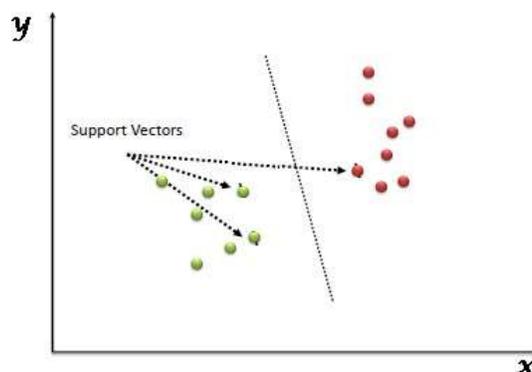


Fig.2 Support Vector Machine

Pseudo code for SVM:

Initialize  $y_i = YI$  for  $i \in I$

REPEAT

Compute SVM solution  $w, b$  for data set with imputed labels

Compute outputs  $f_i = (w, x_i) + b$  for all  $x_i$  in positive bags set  $y_i = \text{sgn}(f_i)$  for every  $i \in I, YI = 1$

FOR (every positive bag  $B_i$ ) IF  $(\sum_{i \in B_i} (1 + y_i)/2 == 0)$

Compute  $i = \text{argmax}_i f_i$  set  $y_i = 1$

End

End

**Pre-processing information**

The coordinated database experienced an information-cleaning process in which we evacuated inappropriate information passages, such as those that gave an unimportant answer, in the database. To smooth woody information, the tuples with an improper passage of information have been wiped out or loaded with the most plausible value, as this is one of the best known techniques to counteract this issue. In addition, the discovery and supplanting work was used to address the irregularity in the organization of the study information.

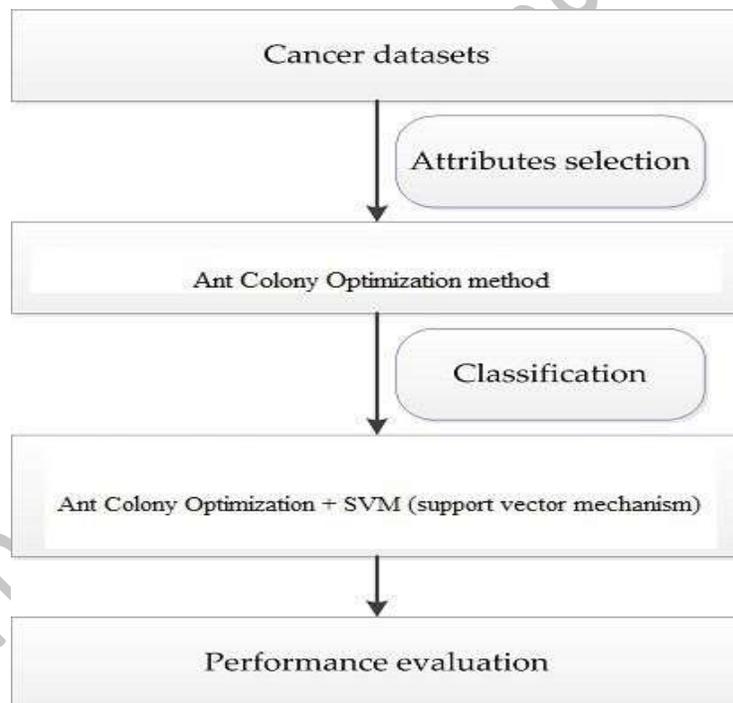


Fig.3 Flow chart to classify cancer data

Figure-3 shows the Case investigation of the classification of cancer information, here at first take the cancer information collection and apply the choice of trait with the assistance of the AOC, and thereafter continue with the AOC-based SVM for the definitive assessment of the evidence and further

apply the cancer expectation to the informational collection.

Feature Selection is the path to finding a subset of features, from the main list of capabilities, in a perfect world. The Insect Colony Optimization System examines how to find a perfect subset of

components using a few cycles. The main objective of the proposed system is to limit the overabundance between them by selecting a subset of features. In this strategy, the most negligible proximity features are chosen by each underground creepy crawl to the features chosen by the past. As such, if the element is chosen by far by most of the ants, this shows that the element has the most negligible equivalence to the exchange of features. The element has the best proportion of pheromone and the chances of securing it through the exchange of ants will be extended in the accompanying cycles. Finally, using the similarity between features, the best features chosen will have a high pheromone effect. The proposed procedure, along these lines, identifies the best features with the least repetition. The element of importance prompts the minimization of the redundancy that will be calculated by taking into account the proximity within the features. The way to highlight ACO 's decision is shown below. In this system, the requested space must be conveyed as an appropriate edge for ACO before the element assurance method begins. Thus, the chase space is conveyed as a totally undirected weighted outline,  $G = \langle F, E \rangle$  where  $F = \{F_1, F_2, \dots, F_n\}$  shows a plan of all features in that each element indicates the center point in the graph,  $E = \{(F_i, F_j) : F_i, F_j \Delta F\}$  shows as far as possible. The relationship between the point of rupture  $(F_i, F_j)$  and  $E$  will be determined by the relationship between  $F_i$  and  $F_j$ .

Stage 1: The explanation behind the proximity of the intersection of features  $I$  and  $j$ . The similarity between any two features is perceived by the calculation of the perfect estimation of the cosine proximity between them. The equivalence of the cosine between the features  $A_n$  and  $B$  is determined using the circumstance.

$$sim(A, B) = \frac{\sum_{i=1}^p I(a_i, b_i)}{\sqrt{\sum_{i=1}^p I^2(a_i)} \sqrt{\sum_{i=1}^p I^2(b_i)}}$$

Here  $A_n$  and  $B$  exhibit two  $p$ -dimensional vector highlights ( $A = \{a_1, a_2, \dots, a_p\}$   $B = \{b_1, b_2, \dots, b_p\}$ ). An estimation of two  $p$ -dimensional vector highlights addressed as 0 and 1 where 1 shows similar highlights and 0 addresses non-practically identical highlights. In the event that the registered similarity of the two highlights is more vital than zero, the

highlights are equal. Using the ACO strategy in the component decision issue, "heuristic data" and "alluring quality" must be represented by the ACO count. In the proposed strategy, heuristic data is depicted as something contrary to closeness within the highlights. An appeal measure of the value of the term "Pheromone" =  $1 \dots n$ , which is related to the highlights and is routinely resurrected by ants.

Stage 2: The proposed system is made up of different emphases. The social pheromone affair allocated to each center point before the cycles begin. In every accent, Nant ants are fixed aimlessly at the various center points. Industriously, as demonstrated by the probabilistic "State Progress Deal," each subterranean insect crosses the centers repeatedly to the point where the cycle stopping principle is satisfied. The end precept is depicted as the number of events in the center points that each insect should select.

$$j = \underset{u \in I}{\text{arg max}} \{[\tau_u][\eta(F_i, F_u)]^q\}, \text{ if } q \leq q_0$$

Stage 3: The State Change Run hopes to pick highlights from most raised pheromones and the smallest similarities to previously selected highlights. Highlight counter show stores the highlights that any insect has picked up. Stage 4: Then, close to the completion of the accentuation, the amount of pheromone for each and every center point is revived by applying the "overall invigorating rule." Considering its component counter, the amount of pheromone for each center point is taken into account. Ants tend to outfit more pheromones to center points with higher counter-element characteristics. What's more, a touch of pheromone breaks down at all center points. Stage 5: The effort is repeated until a given proportion of the cycles is cultivated. Next, the highlights of their pheromone concerns are accumulated in a reduced solicitation. By then, the best-chosen highlights with the most imperative pheromone considerations are chosen as the last subset of the element.

$$P_k(i, j) = \frac{[\tau_j][\eta(F_i, F_j)]^q}{\sum_{u \in I} [\tau_u][\eta(F_i, F_u)]^q}, \text{ if } j \in I; \text{ if } q > q_0$$

Where  $k \ I \ j$  is the unvisited include set,  $\Delta u$  is the pheromone assigned to the element  $u$ ,  $\Delta(F_i, F_u) = 1$ ,

(F Fi u sim is something contrary to the proximity between I and u,  $\beta$  is a parameter used in pheromone versus similarity ( $\beta > 0$ ),  $q_0 [0,1]$  is a consistent parameter, and  $q$  is an unpredictable motivating force in the mean time  $[0,1]$ . In the probabilistic procedure, the expected element  $j$  will be selected from the perspective of the probability  $P_k I j$ ) which is shown as follows:

State progress manages taking into account the parameters  $q$  and  $q_0$ , which are the courses of action between the two.

Operation and Exploration. If  $q \leq q_0$ , then ants choose the perfect element in an eager manner, or there will be outcomes, each component has the credibility of being chosen by its probability, which is determined by the examination. The inspiration behind the probabilistic procedure is to keep away from getting inside a perfect local. The relationship between probabilistic and greedy procedures is

assigned to the "pseudo-subjective comparison rule"<sup>6,7</sup>. The "Overall Revival Guideline" is used for the whole of the center points near the completion of the insect exploration under the following condition:

$$\tau_j(t+1) = (1-q)\tau_j(t) + \frac{FC[j]}{\sum_{j=1}^n 1FC[j]}$$

Where  $n$  is the quantity of one of the species underscores, and the absolute number of pheromone underscores is now and again  $t$  and  $t+1$ , in this manner is the pheromone dissipation  $p$ .

#### 4. Results

The proposed model is implemented using python programming and executed in ANACONDA SPYDER. The proposed model is compared with the traditional models and the results show that the proposed model is better than the traditional methods. The results are depicted below.

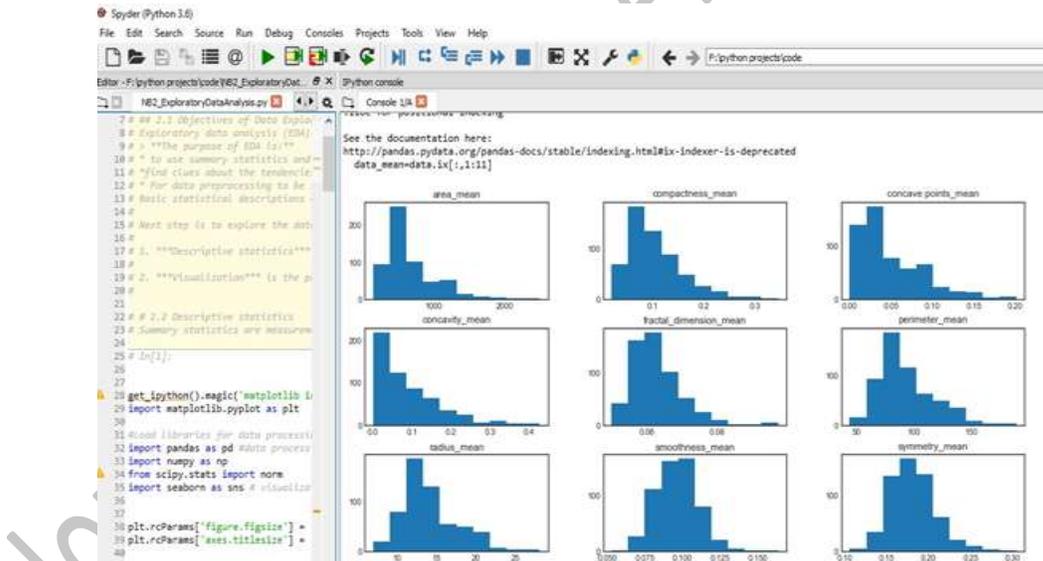


Fig 4 Initial data analysis

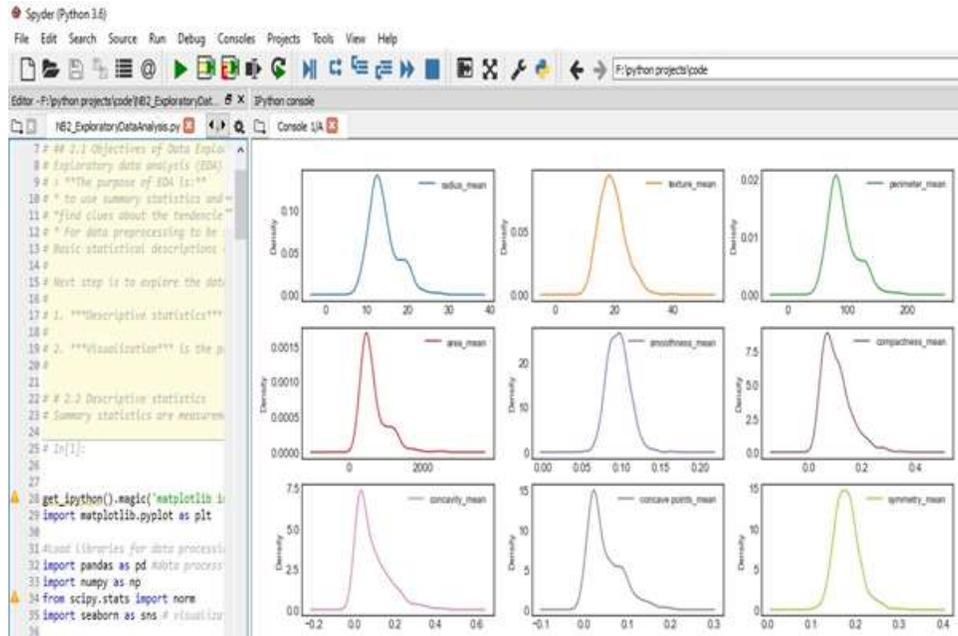


Fig.5 Exploratory data analysis

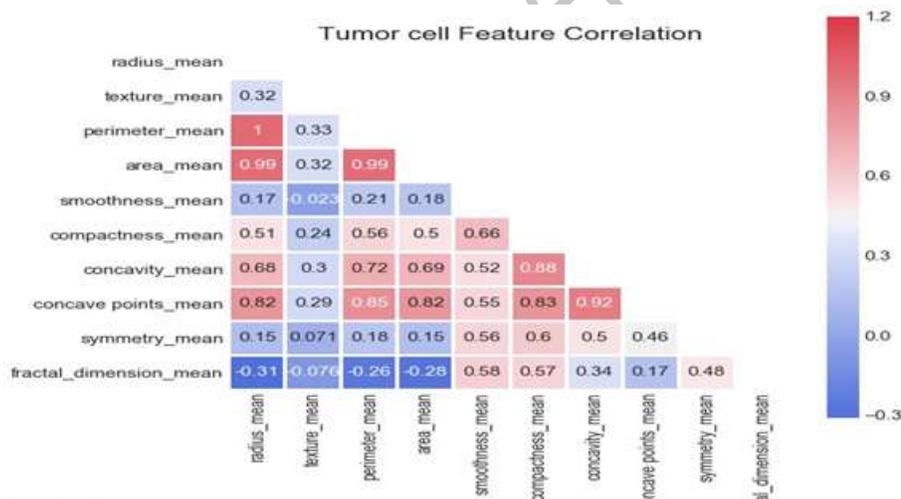


Fig.6 Cancer data features

```
In [3]: runfile('F:/python projects/code/NB3_DataPreprocessing.py', wdir='F:/python
projects/code')
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:444:
DataConversionWarning: Data with input dtype object was converted to float64 by
StandardScaler.
warnings.warn(msg, DataConversionWarning)
```

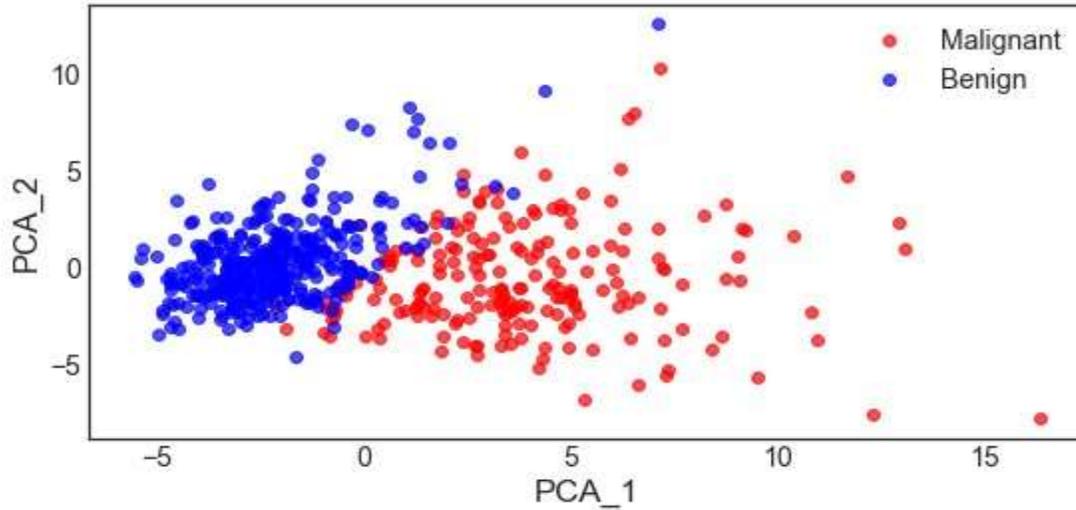


Fig 7 Classification of cancer data with respect to type

The 3-fold cross-validation accuracy score for this classifier is 0.97  
[ 0.93157895 0.95263158 0.94179894]  
Average score and uncertainty: (94.20 +/- 0.496)%

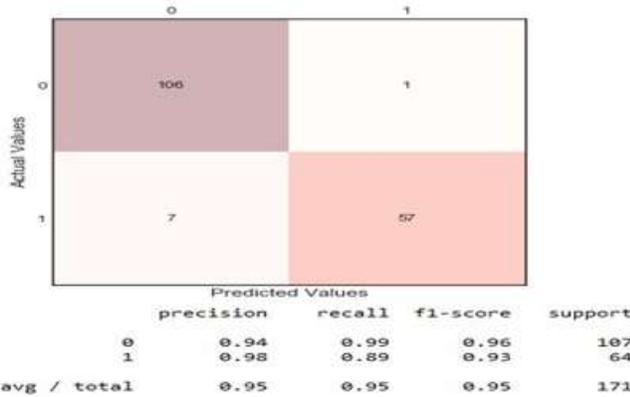


Fig 8 Classification of cancer data with respect to type

### 5. Conclusion

The classification of cancer data is a challenge to data mining. But in this thesis, we work with the SVM-based machine learning combined with the AOC optimization technique. As for accuracy, SVM scores 0.917 in accuracy. In any case, the precision figures

are still higher than that of each existing work, and once again perform best after AOC has been applied, despite the fact that the curacy of air conditioning is declining (0.917). In the context of the other exhibition network, a ton can be resolved with regard to the presentation of the calculations. SVM with AOC scores an ideal 1.000 with respect to the review,

which is crucial as far as the disease expectation is concerned, after AOC has been applied, despite the fact that there are decreases in the estimates of all other presentation measurements of the two calculations mentioned above. Remembering that AOC reduces the exponential to gigantic run time in datasets (both small and huge) and keeping the review score in mind, we can reason that AOC's Logistic Regression and Support Vector Analysis is performing better than Breast Cancer Prediction for this dataset used. As future work can explore more types of deaths, it will be necessary to test and find better experimental results.

### References

- [1]. Hacker Noon Absolute Fundamentals of Machine Learning – Hacker Noon January 15, 2018
- [2]. Furundzic, D.; Djordjevic, M.; Bekic, A.J. Neural networks approach to early breast cancer detection. *J. Syst. Archit.* 1998, 44, 617–633. [CrossRef]
- [3]. Floyd, C.E.; Lo, J.Y.; Yun, A.J.; Sullivan, D.C.; Kornguth, P.J. Prediction of breast cancer malignancy using an artificial neural network. *Cancer* 1994, 74, 2944–2948. [CrossRef]
- [4]. Fogel, D.B.; Wasson, E.C.; Boughton, E.M. Evolving neural networks for detecting breast cancer. *Cancer Lett.* (1995), 96, 49–53. [CrossRef]
- [5]. Fogel, D.B.; Wasson, E.C.; Boughton, E.M.; Porto, V.W.; Angeline, P.J. Linear and neural models for classifying breast masses. *IEEE Trans. Med. Imaging* (1998), 17, 485–488. [CrossRef]
- [6]. Setiono, R. Extracting rules from pruned neural networks for breast cancer diagnosis. *Artif. Intell. Med.* (1996), 8, 37–51. [CrossRef]
- [7]. Wilding, P.; Morgan, M.A.; Grygotis, A.E.; Shoffner, M.A.; Rosato, E.F. Application of backpropagation neural networks to diagnosis of breast and ovarian cancer. *Cancer Lett.* (1994), 77, 145–153. [CrossRef]
- [8]. Wu, Y.; Giger, M.L.; Doi, K.; Vyborny, C.J.; Schmidt, R.A.; Metz, C.E. Artificial neural networks in mammography: Application to decision making in the diagnosis of breast cancer. *Radiology* (1993), 187, 81–87. [CrossRef] [PubMed]
- [9]. H. A. Abbass, “An evolutionary artificial neural networks approach for breast cancer diagnosis.” *Artif. Intell. Med.*, vol. 25, no. 3, pp. 265–81, Jul. 2002.
- [10]. J. Khan, J. S. Wei, M. Ringnér, L. H. Saal, M. Ladanyi, F. Westermann, F. Berthold, M. Schwab, C. R. Antonescu, C. Peterson, and P. S. Meltzer, “Classification and diagnostic prediction of cancers using gene expression profiling and artificial neural networks.” *Nat. Med.*, vol. 7, no. 6, pp. 673–9, Jun. 2001.
- [11]. G.-B. Huang, Q.-Y. Zhu, and C.-K. Siew, “Extreme learning machine: Theory and applications,” *Neurocomputing*, (Dec. 2006.) vol. 70, no. 1–3, pp. 489–501.
- [12]. C. P. Utomo, A. Kardiana, and R. Yuliwulandari, “Breast Cancer Diagnosis using Artificial Neural Networks with Extreme Learning Techniques,” *Int. J. Adv. Res. Artif. Intell.*, vol. 3, no. 7, pp. 10–14, 2014
- [13]. C. P. Utomo, P. S. Pratiwi, A. Kardiana, I. Budi, and H. Suhartanto, “Best-Parameterized Sigmoid ELM for Benign and Malignant Breast Cancer Detection,” pp. 50–55, 2014
- [14]. William H Wolberg, W Nick Street, and Olvi L Mangasarian. (1992). Breast cancer Wisconsin (diagnostic) data set. UCI Machine Learning Repository.
- [15]. Cristianini N, Shawe-taylor J. An Introduction to Support Vector Machines and Other Kernel-based Learning Methods, (2000) London: Cambridge University Press.
- [16]. Joachims T. Making large-scale support vector machine learning practical. *Advances in Kernel Methods: Support Vector Learning.* (1998) MIT Press, Cambridge, MA, 169-184.
- [17]. Machine Learning Mastery Discriminant Analysis for Machine Learning September 22, (2016)

- [18]. DUNTEMAN, G. H. Principal component analysis. quantitative applications in the social sciences series (vol. 69), 1989.
- [19]. Puneet Yadav, RajatVarshney, Vishan Kumar Gupta. Diagnosis of Breast Cancer using Decision Tree Models and SVM (2016)
- [20]. Rohith Gandhi. Nearest Neighbor. Understanding Machine Learning (2018)
- [21]. AdiBronshtein. Train/Test Split and Cross Validation in Python. Understanding Machine Learning (2017).

Journal of Engineering Sciences