

## Lung Cancer Prediction using SVC and SVM Models

Nazma Sultana Shaik<sup>1</sup>, Ummadi Janardhan Reddy<sup>2</sup>, S Nyamathulla<sup>3</sup>, M Kiran Kumar<sup>4</sup>

<sup>1</sup>Research Scholar, <sup>1,2,3</sup>Assistant Professor, <sup>4</sup>Associate Professor

<sup>1,2,3</sup>Department of Information Technology

<sup>1,4</sup>Department of Computer Science and Engineering

<sup>1,2,3</sup>Vignan's Foundation for Science, Technology and Research (Deemed to be University),  
Guntur, Andhra Pradesh, India

<sup>4</sup>Guntur Engineering College, Yanamadala, Guntur, Andhra Pradesh, India

<sup>1</sup>nazma.cs@gmail.com, <sup>2</sup>ummadi.janardan@gmail.com,

<sup>3</sup>nyamath.j@gmail.com, <sup>4</sup>kiran8kumar@gmail.com

### Abstract

Artificial Intelligence domain is the popular in the computer science & it's machine learning application is the core of AI, by these machines will get ability to get learnt in automated manner. This training process will be done through the observation. Computer gives its observation through the data analysis. The Datasets are gathered from different resources like lab tests, patient records from various hospitals & different health survey etc. By this dataset, the evaluation of data point to estimate the results. In the healthcare domain, Machine Learning will be used to predict disease & its risk factors. There are distinct number of algorithms which are used in the ML techniques to make prediction. In this paper, the main crucial approach is to predict the lung cancer using different classification algorithms and comparing those methods.

**Keywords:** Classifiers, K nearest neighbor, Lung Cancer, Machine Learning, SVC, SVM

### I. INTRODUCTION

There has been a gradual advancement of cancer analysis over the past decades. Scientists have used numerous strategies to discover cancer varieties before they cause symptoms, like early stage screening. In addition, they built different models for early estimation of the results of cancer treatment. With the appearance of latest medical technology, giant amounts of cancer knowledge is obtained and square measure accessible to the medical analysis community. Here, one amongst the foremost difficult tasks for physicians is to predict the accuracy of a sickness outcome. As a result, the ml strategies became a preferred medical analysis tool. ml strategies square measure capable of finding and recognizing trends and associations between them from massive datasets, whereas they'll accurately predict future cancer outcomes.

Given the importance of personalised drugs and therefore the increasing trends within the application of ml techniques, we have a tendency to square measure presenting AN analysis of studies exploitation these strategies within the prediction and prognosis of cancer. Prognostic and prognostic characteristics square measure enclosed in these studies. They additionally address the kinds of ml approaches used, the kinds of information that they mix, the general performance of every projected theme.

The utilization of ml approaches would incredibly build the precision of disease risk, recurrence and endurance expectation. With the execution of ML strategies, the precision of malignant growth forecast result has improved impressively by fifteen % – 20 percent as of late.

### II. ML TECHNIQUES

Machine Learning, a branch of Artificial Intelligence, applies the final principle of logical thinking to knowledge sample learning. each learning method consists of 2 phases: [I] estimation of unknown system dependencies from a dataset and [II] use of measured dependencies to predict new system outputs. ml is well-tried to be extremely a very important field of medical specialty analysis with many applications wherever effective generalization is obtained by finding out a

given assortment of biological samples across an n-dimensional space numerous techniques and methodologies.

There are two major groups of common Machine Learning approaches called (I) supervised learning, and (II) unsupervised learning. A labeled assortment of recorded knowledge is employed in supervised learning to approximate or map the input data to the predicted output. Conversely, there are not any tagged examples given to the unsupervised learning methods.

Thusly it is up to the learning plan/model to recognize designs or to find the info information classes. This methodology might be called as challenge of as an order issue in directed learning. The component of arrangement alludes to a learning procedure which classifies the information into a progression of limited classes. Two other regular undertakings including in ML are relapse and grouping. A learning capacity ties the information into a genuine worth vector in the event of relapse issues.

The prognostic variable value will then be determined for every new sample supported this approach. Clustering could be a common unattended activity with in which one makes an attempt to search out the classes or clusters to classify the things of the information. On the idea of this method, every new sample could also be assigned to 1 of the clusters known for the similar characteristics they share.

In general, a classification algorithmic rule could be a methodology that weights the input characteristics so the output divides one class as positive values, in addition as into negative values. Classifier coaching is performed to work out the weights (and features) which give the foremost correct and best differentiation between the 2 knowledge categories. Linear discriminant analysis is that the best classifier to classify the linear coefficient of complex outcomes as some way of optimizing the discrepancy between the 2 groups suggests that. However, the relative distinction between the 2 teams is not well drawn by one line for several knowledge sets.

Supporting vector machines, artificial neural networks, and random forest trees are measure newer machine strategies that may produce a lot of advanced distinctions between the 2 knowledge categories, and every of those classifiers has benefits and drawbacks.

### III. LITERATURE SURVEY

ML has been used in various domains, in healthcare it plays a major role in early prediction/diagnosis. Cancer is one of the primary cause of death, to predict in early stage ML has become a good tool for researchers. ML works by identifying pattern [2], clustering it and making prediction based on predicted outcomes. Number of classification algorithms are implemented to find the accuracy of algorithms [3].

#### [1] Classification of Cancerous Profiles using Machine Learning:

The author used gene expression to identify drug treatment, because recommended drug for individual varies by following factor [i] cancer type [ii] severity [iii] genetic heterogeneity which is the most important factor. The proposed model used hybrid algorithm that contain both inner and outer classification. The algorithm's are Clustering using Neural Network and Classification using SVM. The result has been evaluated using Precision, Re-call, F-Measure, Accuracy. The accuracy of algorithm's SVM-93.57% and KNN-87.51.

#### [4] Data Processing Techniques in Multiple Cancer Prediction:

The main objective of the paper is to demonstrate how Data Mining techniques like Association mining, Clustering and Classification are used for lung cancer prediction.

#### [5] Application Of K-Nearest Neighbor (K-NN) Algorithm On Lung Disease Diagnosis Expert System:

The focus of study is on application of KNN algorithm on diagnosis of different types of

lung disease, with a aim to design application which inherits KNN algorithm on expert system .The result of prediction using KNN algorithm is 96.97%.

#### IV. PROPOSED SYSTEM

The technique proposed starts with the collection of data which is accompanied by preprocessing. The selected classifiers will then be trained and checked on the dataset for benchmark. The tests are measured and evaluated to determine the best methodology for detecting lung cancer.

##### A. K Nearest Neighbour:

K-Nearest Neighbour algorithm is a simple classification strategies of system study in set of rules. It has been used in many device learning application. By this algorithm, the classifying objects assigns the reference space characteristics vector (training statistics) to the label of its K-nearest neighbour or neighbours. It is defined by using distance and quantity of K, so this set of rules is known as closest/nearest neighbouring set of rules. The KNN algorithm is one of the exceptional and is extensively used for classification and clustering. In this process, each pattern cost fits into the test dataset and is graded on the premise of the training facts consistent with the nearest k samples. Values of the class numbers derived from sample value of k, the most number is calculated from samples of class. KNN is a easy clustering set of rules that classifies information set based totally on their neighbourly similarity. K stands for variety of gadgets set for the type which might be considered. The check pattern mark is determined via the closest suit the various nearest k neighbour. Euclidean Distance (ED) is utilized in this approach to discover the distance among the test samples and study of samples.

##### B. Support Vector Classifier:

A Linear SVC (Support Vector Classifier) makes an attempt to suit the info returning a "best fit" hyperplane that separates the info or classify it. (fig 4.2.1). Then we will feed some options to our classifier once obtaining the hyperplane, to visualize what the "predicted" category is. This makes this specific formula terribly ideal for our desiers, though this could be employed in alternative cases

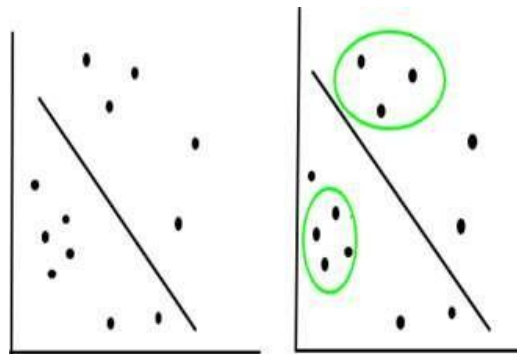
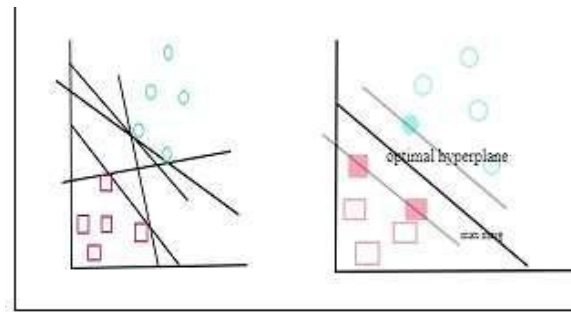


Fig 4.2.1

##### C. Support Vector Machine:

SVMs are a more modern approach to cancer prediction/prognosis approaches enforced in ml. SVMs map the input vector into a higher-dimensionality feature space and outline the hyperplane (fig four.3.1) that divides the info points into 2 groups [10]. This maximizes the marginal distinction between the choice hyperplane and also the patterns nearest to the boundary [11]. The ensuing classifier achieves tremendous generalizability, therefore it will be wont to accurately determine new samples. it's value noting that one may acquire probabilistic outputs for SVMs [12-15].

**Fig 4.3.1**



## V. TECHNOLOGY USED ANACONDA

Anaconda distribution could be a free and open- source ,used to assign these packages for Python and R programming languages for clinical computing such as information science,gadget gaining knowledge of packages,large-scale record processing ,predictive analytics,etc.,with the goal of remodelling package control and planning.Anaconda has several more packages as conda application and virtual atmospheres,It also include the user interface called the Anaconda Navigator,as a graphical interface become independent from the declaration interface.

### ANACONDA NAVIGATOR

This guide may be a work area graphical interface encased in boa constrictor pilot dissemination that permits the client to dispatch applications and oversee conda bundles, air and channels. Guide will scavenge around for bundles on Anaconda pilot Cloud or in a local Anaconda pilot Repository, setup them in an air, run the bundles and overhaul them. It's available for Windows, macOS and UNIX.

The following applications are on the market by default in Navigator

- JupyterLab
- Jupyter Notebook
- Spyder

### SPYDER

Spyder (Scientific Python Development Environment) is one among open flexibly cross-stage incorporated improvement air (IDE) for logical programming inside the Python language. Spyder incorporates with number of amazing bundles inside the logical Python stack, just as NumPy, SciPy, Matplotlib, pandas, IPython, SymPy and Cython, further as elective open source software. It's available cross-stage through Anconda.

### PYTHON

Python is an interpretive, high-level, widespread programming language. The goal of language is to assist programmers write simple, logical, small- large- code.It supports multiple programming paradigms in addition to procedural,objective and practical programming.

## VI. RESULT AND DISCUSSION

The medical data are used to find the accuracy of different algorithm. Here we compared the accuracy of three different algorithms namely, KNN, support vector classifier, support vector machine. The accuracy is also predicted using confusion matrix.

Table-I: Accuracy of algorithm's

Algorithm used	Medical Dataset
KNN	93.1%
SVC	86.2%
SVM	100%

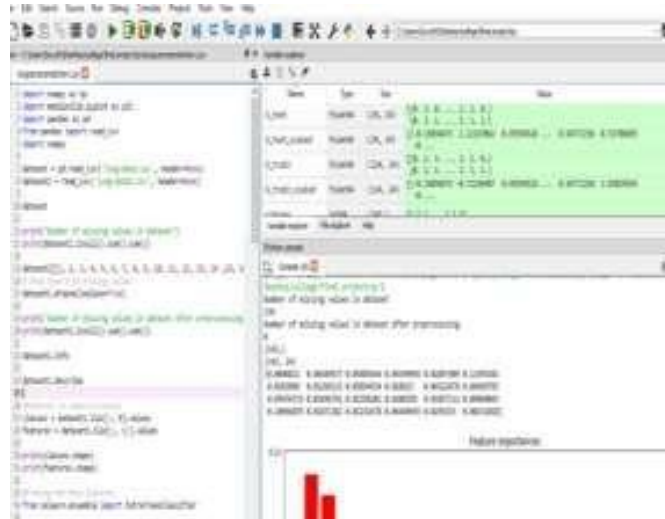


Fig 6.1

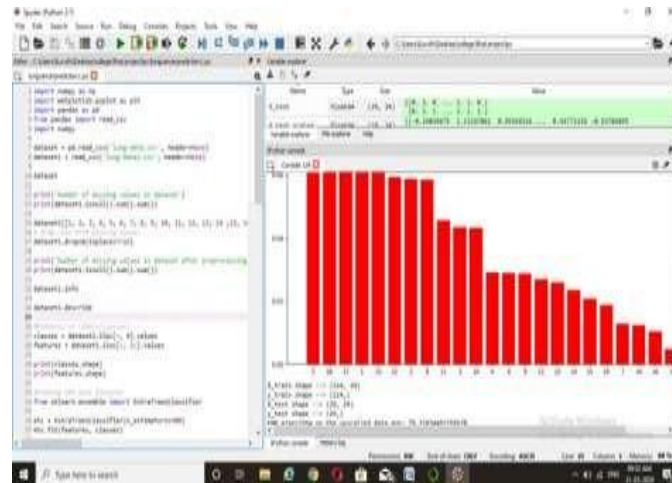


Fig 6.2

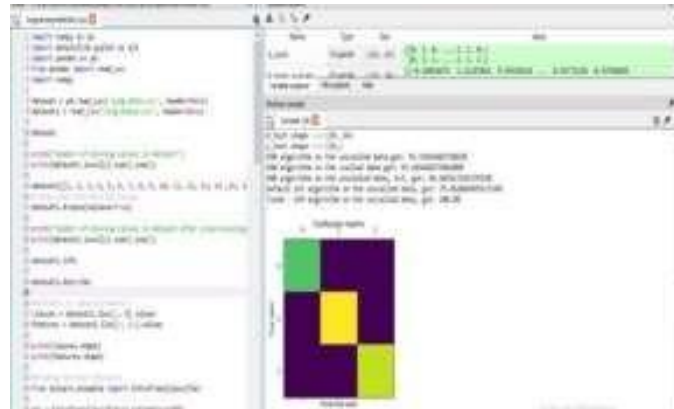


Fig 4.3

## VII. CONCLUSION

Here, we investigated the ideas of ML in this report, along these lines plotting the ml application in disease expectation/visualization. A considerable lot of the investigations proposed are focused on the advancement of gauging models utilizing directed ML techniques and grouping calculation to anticipate genuine infection results. It is clear from the examination of their discoveries that the reconciliation of confounded heterogeneous information, joined with the use of different procedures for the choice and arrangement of highlights, can give promising induction instruments in the field of malignant growth.

## REFERENCES:

1. Aman Sharma, Rinkle Rani Department of computer Sc. & Engg “Classification of Cancerous Profiles mistreatment using Machine Learning.” International Conference on Machine learning and knowledge Science, 2017
2. P. Ramachandran Ph.D Research Scholar, Department of CS&A, N. Girija, Ph.D Lecturer, Department of IT, T. Bhuvaneshwari, Ph.D Asst. Professor, Department of CS&A, “Early Detection and Prevention of Cancer cells using Data Processing Techniques”. 2014
3. Muhammad Imran , Saba Bashir, Zain Khan, Farhan Hassan Khan Department of Computer Science “An Analysis of Machine Learning Classifiers and Ensembles for Early Stage Prediction of Lung Cancer”. 2018
4. Dr. A. R. Pon Periasamy , K. Arutchelvan “ Data Mining Techniques in Multiple Cancer Prediction” ISSN: 2277 128X Volume 7, May 2017.
5. Olha Musa, A. Malik Buna, Marlin, R. Rizal Isnanto, Suryono Academic of Management and Informatica of PC “Application Of K-Nearest Neighbor (K-NN) Algorithm rule On respiratory Organ Disease Diagnosis Expert System” International Journal of Scientific & Engineering Research Volume nine, Issue 12, December-2018.
6. R. Sathishkumar, K. Kalaiarasan, A. Prabhakaran, M. Aravind, dept of computer science “Detection of lung cancer using SVM classifier and KNN algorithm” International journal of scientific research and review Volume 8, Issue 3, 2019.
7. N. Sudhir, V. Khanaa, “Detection and Prediction of Lung Cancer Using Various Algorithms” International Journal of Engineering and Advanced Technology (IJEAT) ISSN: 2249-8958, Volume-8, Issue-6S3, September 2019.
8. Mr. Sandeep A. Dwivedi 1 , Mr. R. P. Borse 1 , Mr. Anil M. Yametkar 2 1 Department of E & TC, SAE, “Lung Cancer detection and Classification by using Machine Learning & Multinomial Bayesian”, IOSR Journal of Electronics and Communication Engineering (IOSR-JECE) e-ISSN: 2278-2834, p-ISSN: 2278- 8735. Volume 9..
9. Zehra , Taner Tunç 2 Ondokuz Mayıs University Samsun/Turkey, “Lung Cancer Detection and Classification with Classification Algorithms”, IOSR Journal of Computer Engineering (IOSR-JCE) e-ISSN: 2278-0661, p-ISSN: 2278-8727, Volume 18, Issue 6, Ver. III (Nov.-Dec. 2016), PP 71-7.
10. G. Ramu, P. Dileep Kumar Reddy, Appawala Jayanthi “A Survey of Precision Medicine Strategy

Using Cognitive Computing” International Journal of Machine Learning and Computing, Vol. 8, No. 6, December 2018 DOI: 10.18178/IJMLC2018.8.6.741 (Scopus) (UGC Approved) Journal No: 48748, pp 530 to 535.

11. Ummadi Janardhan Reddy, Pandluri Dhanalakshmi, Pallela Dileep Kumar Reddy Image Segmentation Technique Using SVM Classifier for Detection of Medical Disorders Ingénierie des Systèmes d’Information, Vol. 24, No. 2, pp. 173-176, April 2019, <https://doi.org/10.18280/isi.240207> (Scopus) ISSN: 1633-1311 (print); 2116-7125 (online) Impact Factor : 0.409
12. Singamaneni Kranthi Kumar, Pallela Dileep Kumar Reddy, Ramesh G, Venkata Rao Maddumala (2019). Image transformation technique using steganography methods using LWT technique. Traitement du Signal, Vol. 36, No. 3, pp. 233-237. June 2019, <https://doi.org/10.18280/ts.360305>, (WOS, SCI- E) (UGC Care List) ISSN: 0765-0019 (print); 1958-5608 (online) Impact Factor : 0.387.
13. J. Somasekar a, , G. Ramesh , Gandikota Ramu, P. Dileep Kumar Reddy, B. Eswara Reddy e, Ching-Hao Lai, “A dataset for automatic contrast enhancement of microscopic malaria infected blood RGB images”, Data in brief, Elsevier, <https://doi.org/10.1016/j.dib.2019.104643>,2352-3409/2019
14. Ramu, G., Jayanthi, A, “A review on precision medicine and its advantages”, Pakistan Journal of Biotechnology 14(3):515-519 · January 2017
15. Ramu, G., Soumya, M., Jayanthi, A. et al., “Protecting big data mining association rules using fuzzy system”, TELKOMNIKA, Vol.17, No.6, December 2019, pp.101~109, ISSN: 1693-6930,