

LUNG CANCER DETECTION USING SUPPORT VECTOR MACHINE

N.Ashok Kumar, M.Vyshali, P.Sravya, J.Dedipya, S. Alekhy

Department of Computer Science and Engineering, Vignan's Nirula Institute of Technology and Science for Women, Pedapalaluru, AP, India.

Abstract

The Lung Cancer detection is a very challenging task. It generally occurs in both male and female due to uncontrollable growth of cells in lungs. Thus, treatment quality is a key stage to improve the quality of life of patients. Machine Learning has a great influence to healthcare sector because of its high computational capability for detection of disease with accurate data analysis. Hence, the trusted and automatic classification scheme is essential to prevent the death rate of humans. Support Vector Machine (SVM) algorithm in Machine Learning (ML) is used to detect the Lung Cancer by classification. The Lung Cancer detection is divided into two phases such as training and testing phases. In the training phase, preprocessing, feature extraction and SVM classification is performed to make a prediction model more efficient. Finally, Support Vector Machine is used for Lung Cancer detection. The cancer dataset is taken from Kaggle. The proposed method detects the lung cancer with best accuracy.

1. Introduction:

Machine Learning algorithms have the potential to be invested deeply in all the fields of medicine from drug discovery to clinical decision making, significantly altering the way medicine is practiced. The success of machine learning algorithms at computer vision tasks in recent years comes at an opportune time when the medical records are increasingly digitalized. Medical images are an integral part of a patient's health and are currently analysed by human radiologists, who are limited by speed, fatigue, and experience. And it also takes years, great financial cost to train a qualified radiologist. A delayed or erroneous diagnosis may cause harm to the patients. Therefore, it is an ideal for medical image analysis to be carried out by an automated, accurate and efficient machine learning algorithms. Lung cancer is the most common and aggressive disease, leading to a very short life expectancy in their highest grade. The major types of lung cancer are non-small cell lung cancer and small cell lung cancer. They differ in size of the cell. Reports say around 84% of lung cancer cases are non-small cell and 13% are small cell lung cancer. Diagnosis of lung cancer is done medically. In general, the patient will be aware when the cancer is in advanced stage. Accidentally sometimes patients may get the existence of malignant cells early and if a person has a cough and generating sputum, testing the sputum beneath the microscope can periodically reveal the lung cancer cells. When moving images are considered, MRI scan does not work well like in case of lungs which expands as we breathe in and contracts when we breathe out. In these severe cases it can affect brain, bones or distant sites. The major focus is to predict the risk level of lung cancer. In the past, the outcome for patients diagnosed with these tumours was very poor, with typical survival rates of just several weeks. More sophisticated diagnostic tools, in addition to innovative surgical and radiation approaches, have helped survival rates expand up to years and also allowed for an improved quality of life of patients. The classification techniques helped to detect the survival possibilities of lung cancer victims and help the physician to take decision on the forecast of disease.

2. Literature survey:

All researchers have aimed to develop a system which will predict and detect the cancer in its early stages. Also tried

to improve the accuracy of the Early Prediction and Detection system by pre-processing, feature extraction and classification techniques of extracted database.

S Vishukumar, K. Patela and Pavan Shrivastava - In this paper authors mostly focus on significant improvement in contrast of masses along with the suppression of background tissues is obtained by tuning the parameters of the proposed transformation function in the specified range. The manual analysis of the sputum samples is time-consuming, inaccurate and requires intensive trained persons to avoid diagnostic errors. The segmentation results will be used as a base for a Computer Aided Diagnosis (CAD) system for early detection of cancer, which improves the chances of survival for the patient. In this paper, authors proposed Gabor filter for enhancement of medical images. It is a very good enhancement tool for medical images.

3. Existing System:

Random Forest Classifier is most accurate algorithm mainly used in healthcare domain. This algorithm is applied on data set to achieve accuracy which helps in predicting lung cancer at early stage. Random forest builds multiple decision trees and merges them together to get accurate and stable prediction. It will handle the missing values and maintain the accuracy of a large proportion of data.

Disadvantages of Existing System:

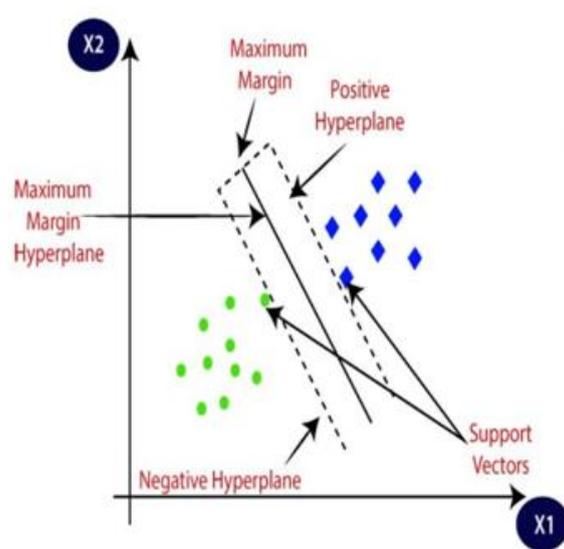
- Random Forest requires more time to train as it generates a lot of trees.
- Requires more computational power and resources.
- Makes decision on majority of voting.
- Low accuracy with more complexity.

4. Proposed System:

Support Vector Machine is powerful and flexible supervised machine learning algorithm used for classification and is popular because of its ability to handle multiple and categorical attributes.

Support Vector Machine algorithm is implemented with kernel that transforms an input data space into required form. Kernel makes

SVM more powerful, accurate and flexible. In our project SVM algorithm is implemented by using Linear Kernel. It can be used when there are a large number of features in particular dataset. Training a linear kernel is faster than other kernels. The objective of a Linear SVC is to fit to the data we provide and returning a best fit hyperplane that divides or categorizes the data.



There are some important concepts in SVM like Hyperplane, support vectors, Margin. Hyperplane is used to divide a set of objects having different classes. Support vectors are the data points closest to hyperplane. Margin is the gap between two lines on closest data points of different classes. At first, SVM will generate hyperplanes iteratively which segregates the classes in the best way. Then it will choose the best hyperplane that separates the classes correctly. It uses a technique called the kernel trick to transform the data and based on these transformations it finds an optimal boundary between possible outputs.

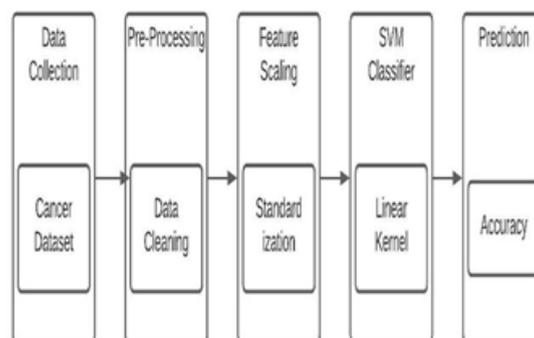


Fig: Process Flow Diagram

The above figure shows the proposed model. First, the cancer dataset is collected from Kaggle. At the Pre-Processing phase, data cleaning is performed which deletes unwanted columns like Patient ID and then Feature Scaling is done. It is a technique which will standardize the independent features present in the data in a fixed range. It re-scales a feature value so that it has a distribution with 0 mean value and variance equals to 1. SVM classifier uses training data set and generates classification rules. By using these classification rules with test data set SVM classifier predicts the output.

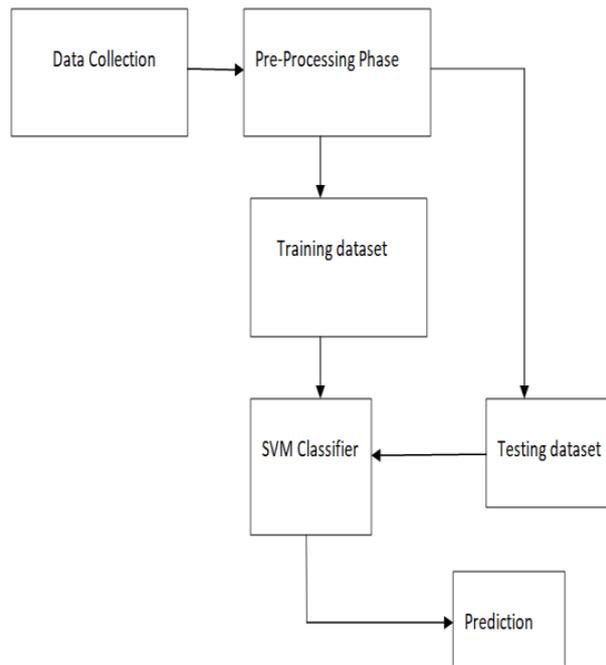


Fig: System Architecture

5. Result:

The Lung Cancer detection is divided into two phases such as training phase and testing phases. The dataset is collected from Kaggle which consists of 1300 records. The training dataset consists of 1000 records and testing dataset consists of 300 records. By using Support Vector Machine, the classification report is generated with high accuracy of 98.33%. The SVM algorithm classifies the records into different classes based on the cancer level of patient. From the below output we can say that 89 people fall under class-0 with low cancer, 96 fall under class-1 with medium cancer and 115 fall under class-2 with high cancer.

```

**
Accuracy Using SVM: 0.9833333333333333
      precision    recall  f1-score   support

   0       0.95      1.00      0.97         89
   1       1.00      0.95      0.97         96
   2       1.00      1.00      1.00        115

 accuracy          0.98         300
 macro avg         0.98         300
 weighted avg      0.98         300
  
```

6. Conclusion:

The main goal is to design an efficient model for classification of lung cancer with high accuracy, performance and low complexity. The random forest classification is performed by using pre-processing techniques and the complexity is low but the computation power is

high mean while accuracy is also low. Further to improve the accuracy and to reduce the computation time, Support Vector Machine (SVM) classification algorithm is introduced in the proposed system. Also, the classification report is given with the high accuracy of 98.33%. On the basis of this information the best therapy, surgery, radiation and chemotherapy is advised. So, early detection is needed to save the lives from the disease.

7. Future Scope:

It is a challenging task in machine learning to construct a specific and computationally efficient classifier for the medical applications. Future scope of this project is how the big datasets will behave for the classification algorithms and the identification of particular stage of lung cancer is done in near future. To identify the particular stage, we need to perform the detection on image dataset using Deep Convolutional Neural Networks. A few innovative applications that span across traditional medical image analysis categories are described here. The first one is the lung cancer detection which will classify the lung CT images into benign or malicious and thus form the dataset abnormal lung CT can be extracted. At last, after segmentation lung cancer stage identification is done which will find out the stages of all nodes that have been segmented. Based on the size of nodules, it is classified into different stages which will help to diagnose patient's degree of spread of disease more accurately than the previous existing methods.

References

- [1] Kwetishe Joro Danjuma, "Performance Evaluation of Machine Learning Algorithms in Post-operative Life Expectancy in the Lung Cancer Patients" Department of Computer Science, Modibbo Adama University of Technology, Yola, Adamawa State, Nigeria.
- [2] Survey of Intelligent Methods for Brain Tumor Detection-IJCSI International Journal of Computer Science Issues, Vol. 11, Issue 5, No1, September 2014.
- [3] Zehra Karhan¹, Taner Tunç², "Lung Cancer Detection and Classification with Classification Algorithms" IOSR Journal of Computer Engineering (IOSR-JCE) e-ISSN: 2278-0661, p-ISSN: 22788727, Volume 18, Issue 6, Ver. III (Nov.-Dec. 2016), PP 71-77.
- [4] Ada, Rajneet Kaur, "A Study of Detection of Lung Cancer Using Data Mining Classification Techniques" International Journal of Advanced Research in Computer Science and Software Engineering, Volume 3, Issue 3, March 2013.
- [5] Survey of Intelligent Methods for Brain Tumor Detection-IJCSI International Journal of Computer Science Issues, Vol. 11, Issue 5, No1, September 2014.
- [6] 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, Issue 1, Ver. III (Jan. 2014), PP 69-75.
- [7] K.V. Bawane, A.V. Shinde "Diagnosis Support System for Lung Cancer Detection Using Artificial Intelligence"- International Journal of Innovative Research in Computer and Communication Engineering, Vol. 6, Issue 1, January 2018.
- [8] H.R. HAl-Absi, B.B. Samir, K. B. Shaban and S. Sulaiman, "Computer aided diagnosis system based on machine learning techniques for lung cancer", 2012 International Conference on Computer and Information Science (ICCIS), Kuala Lumpur, 2012, pp. 295-300.
- [9] Sukhjinder Kaur "Comparative Study Review on Lung Cancer Detection Using Neural Network and Clustering Algorithm", International Journal of Advanced Research in Electronics and Communication Engineering (IJARECE) Volume 4, Issue 2, February 2015.
- [10] D. Vinitha, Dr. Deepa Gupta, and Khare, S., "Exploration of Machine Learning Techniques for Cardiovascular Disease", Applied Medical Informatics, vol. 36, pp. 23-32, 2015.
- [11] Sathyan Panicker, J.V., "Lung Nodule Classification Using Deep Conv Net on CT Images", 9th International Conferen

ceonComputing,CommunicationandNetworking Technologies, ICCCNT2018.

- [12] Isaac,J.,Harikumar,“LogisticregressionwithinDBMS”
,Proceedingsofthe20162ndInternationalConferenceonContemporary Computing.
- [13] PengGuan,DeshengHuang,Miao He, and Baosen Zhou,
“LungCancerGeneExpressionDatabaseAnalysisIncorporatingPriorKnowledgewithSupportVectorMachineBase
dClassificationMethod”,JExpClinCancerRes,28(1): 103,2009.
- [14] A. K. Santra, and C. JosephineChristy,“GeneticAlgorithmandConfusionMatrixforDocumentClustering”,IJCSIIInternationalJournal
of Computer Science Issues,Vol.9,Issue1,No.2,2012.