

International Journal of
Engineering Research and Science & Technology



ISSN : 2319-5991

www.ijerst.com

Email: editor@ijerst.com or editor.ijerst@gmail.com

LUNG CANCER DETECTION USING MACHINE LEARNING

1 DR. ABDUL RAHIM, 2 P. SAHITHI, 3 R. RUCHITHA, 4 U. SAMATHA, 5 D. SOUMYALATHA

ABSTRACT: The Main Objective of this research paper is to find out the early stage of lung cancer and explore the accuracy levels of various machine learning algorithms. After a systematic literature study, we found out that some classifiers have low accuracy and some are higher accuracy but difficult to reach nearer of 100%. Low accuracy and high implementation cost due to improper dealing with DICOM images. For medical image processing many different types of images are used but Computer Tomography (CT) scans are generally preferred because of less noise. Deep learning is proven to be the best method for medical image processing, lung nodule detection and classification, feature extraction and lung cancer stage prediction. In the first stage of this system used image processing techniques to extract lung regions. The segmentation is done using K Means. The features are extracted from the segmented images and the classification are done using various machine learning algorithm. The performances of the proposed approaches are evaluated based on their accuracy, sensitivity, specificity and classification time.

Keywords: Structural Co-occurrence Matrix (SCM), Classifier, Data Set, ROC curve, Malignant nodule, Benign nodule

INTRODUCTION

The cause of lung cancer stays obscure and prevention become impossible hence the early detection of lung cancer is the only one way to cure. Size of tumour and how fast it spread determine the stage of cancer [1]. Lung cancer spreading widely all over the world. Death and health issue in many countries with a 5-year survival rate of only 10–16% [2][3]. In some cases, the nodules are not clear and required a trained eye and considerable amount of time to

detect. Additionally, most pulmonary nodules are not cancerous as they can also be due to non-cancerous growths, scar tissue, or infections [4]. Even though many researchers use machine learning frameworks. The problem with these methods is that, in order to evaluate the best performance, many parameters need to be hand-crafted which is making it difficult to reproduce the better results [5].

1 ASSISTANT PROFESSOR, DEPARTMENT OF ECE, MALLA REDDY ENGINEERING COLLEGE FOR WOMEN, HYDERABAD
2,3,4&5 UG SCHOLAR, DEPARTMENT OF ECE, MALLA REDDY ENGINEERING COLLEGE FOR WOMEN, HYDERABAD

Classification is an important part of computation that sort images into groups according to their similarities [6][7]. In the structure of cancer cell, where most of the cells are overlapped with each other. Hence early detection of cancer is more challenging task [8][9]. After an extensive study, we found that ensemble classifier was performed well when compared with the other machine learning algorithms [10]. The existing CAD system used for early detection of lung cancer with the help of CT images has been unsatisfactory because of its low sensitivity and high False Positive Rates (FPR).

LITERATURE REVIEW In paper [11] Pankaj Nanglia, Sumit Kumar et al proposed a unique hybrid algorithm called as Kernel Attribute Selected Classifier in which they integrate SVM with Feed-Forward Back Propagation Neural Network, which helps in reducing the computation complexity of the classification. For the classification they proposed three block mechanisms, pre-process the dataset is the first block. Extract the feature via SURF technique followed by optimization using genetic algorithm is the second block and the third block is classification via FFBPNN. The overall accuracy of the proposed algorithm is 98.08%. In paper [12] Chao Zhang, Xing Sun, Kang Dang et al perform a sensitivity analysis using the multicenter data set. They chosen two categories Diameter and Pathological result. Diameter were divided into three sub groups. 0-10mm, 10-20mm, 20- 30mm. In 0-10mm group sensitivity 85.7% (95% CI, 70.8%-100.0%) and specificity 91.1% (95% CI, 86.8%-95.2%) were found. In 10-20mm group sensitivity 85.7% (95% CI, 77.1%-94.3%) and specificity 90.1% (95% CI, 84.8%-95.4%) were found. In 20-30mm group sensitivity 78.9% (95% CI, 66.0%-91.8%) and specificity 91.3% (95% CI, 83.2%-99.4%) were found. The algorithm had provided the highest accuracy of 85.7% for adenocarcinoma and 65.0% for Squamous cell carcinoma.

In paper [13] Nidhi S. Nadkarni and Prof. Sangam Borkar focuses their study mainly on the classification of lung images as normal and abnormal. In their proposed method median filter was used to eliminate impulse noise from the images. Mathematical morphological operation enables accurate lung segmentation and detect tumour region. Three geometrical features i.e. Area, perimeter, eccentricity was extracted from segmented region and fed to the SVM classifier for classification. In paper [14] Ruchita Tekade, Prof. DR. K. Rajeswari studied the concept of lung nodule detection and malignancy level prediction using lung CT scan images. This experiment has conducted using LIDC_IDRI, LUNA16 and Data Science Bowl 2017 datasets on CUDA enabled GPU Tesla K20. The Artificial Neural Network used to analyze the dataset, extracting feature and classification purpose. They used U-NET architecture for segmentation of lung nodule from lung CT scan images and 3D multigraph VGG like architecture for classifying lung nodule and predict malignancy level. Combining these two approaches have given the better results. This approach given the accuracy as 95.66% and loss 0.09 and dice coefficient of 90% and for predicting log loss is 38%. In paper [15] Moffy Vas, Amita Dessai, studied mainly on the classification of lung images cancerous and non-cancerous. In their proposed method pre-processing was done, in which unwanted portion of the lung CT scan was removed. They used median filter to eliminate salt and pepper noise. Mathematical morphological operation enables accurate lung segmentation and detect tumour region. Seven extracted features i.e. energy, correlation, variance, homogeneity, difference entropy, information measure of correlation and contrast respectively was extracted from segmented region and fed to the feed forward neural network with back propagation algorithm for classification. The algorithm looks for the least of the

error function in the weight space gradient descent method. The weights are shuffled to minimise the error function. The training accuracy was 96% and testing accuracy was 92%. The sensitivity was 88.7% and specificity was 97.1%. In paper [16] Radhika P R, Rakhi.A.S.Nair, mainly focused on prediction and classification of medical imaging data. They used UCI Machine Learning Repository and data.world. dataset. Used various machine learning algorithm for comparative study and found that support vector machine gives higher accuracy 99.2%. Decision Tree provide 90%, Naïve Bayes provide 87.87% and Logistic Regression provide 66.7%. In paper [17] Vaishnavi. D1, Arya. K. S2, Devi Abirami. T3 , M. N. Kavitha4, studied on lung cancer detection algorithm. In pre-processing they used Dual-tree complex wavelet transform (DTCWT) in which the wavelet is discretely sampled. GLCM is second order statistical method for texture analysis which provide a tabulation of how different combination of Gray level co-occur in an image. It measures the variation in intensity at the pixel of interest. They used Probability Neural Network (PNN) classifier evaluated in term of training performance and classification accuracy. It gives fast and accurate classification. In paper [18] K.Mohanambal , Y.Nirosha et al studied structural co-occurrence matrix (SCM) to extract the feature from the images and based on these features categorized them into malignant or benign. The SVM classifier is used to classify the lung nodule according to their malignancy level (1 to 5).

SYSTEM MODEL

A. DATA EXPLORATION Three datasets are used in this research containing labelled nodules positions for image segmentation and cancer/non-cancer labels for classification [19].

1. TCIA Dataset The cancer imaging archive (TCIA) host collection of de-identified medical images, primarily in DICOM format. Collections are organized

according to disease and image modality (such as MRI or CT). CT images data used to support the findings of this study have been deposited in the Lung CT-Diagnosis repository

2. Lung Image Database Consortium Image Collection (LIDC-IDRI) consists of lung CT scans of 1018 patients (124GB) in DICOM format. Four experienced radiologists independently reviewed the lung CT scans and annotated the nodules in the dataset. 3. Kaggle data science bowl 2017 provides lung CT scans of 1595 patients (146GB) in DICOM format and having a set of labels, which denote that if the patient was diagnosed with lung cancer in future, even one year after the scan were taken.

B. ALGORITHMS AND TECHNIQUES

The U-Net Convolutional Network is used for biomedical image segmentation. It takes an input image and an output mask of the region of interest. It first generates a vector of features typically in a convolutional neural network, and then use another upconvolutional neural network to predict the mask given by the vector of features [20][21][22]. This is a binary classification task using morphological and radiological features extracted from the images and masks. The features are continuous and numerical, but can be discretized into categories. The following classifiers were explored [23][24][25].

1. Logistic regression is particularly strong in binary classification which provide top candidate model for completion of this task.

2. Gaussian Naïve Bayes is suitable for the continuous numerical features. It takes the mean and variance for each feature in each class [26].

3. Multinomial Naïve Bayes required the categorical data. In this feature transformed into discrete steps. This may be more suited than Gaussian NB since some of the feature distributions representing a class is not normally distributed. For example, diameter with

non-cancer is strongly skewed to the left [27].

4. Support Vector Machines draws a separation line that maximizes the points representing the classes in a multidimensional feature space. A kernel trick can be used to fit a more defined boundary [28].

5. Random Forest frequently used on kaggle for classification tasks. It creates many decisions trees with random samples and features and takes a vote on its output. This is used to prevent overfitting.

6. Gradient Boosting also frequently used on kaggle for classification tasks. It's similar to Random Forest but instead of random samples for each tree, it takes the samples with the highest error on the previous tree to train the successive trees.

7. Ensemble classifiers are created by averaging the output of several of the above models.

C. MODEL EVALUATION AND VALIDATION Model 1: U-Net Convolutional Neural Network for nodule segmentation [29].

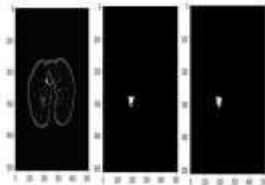


Figure 1 U-Net image segmentation. Processed CT image (left), ground truth label (center), predicted label (right)

CONCLUSION CAD system for lung cancer includes the stages of pre-processing, nodule detection, nodule segmentation, feature extraction and classification of the nodule as benign or malignant. Once the nodules are detected and segmented the feature extraction process begins. The features necessary for classification are extracted using feature extraction techniques from the segmented nodule. Based on the features extracted, a classifier is used for classifying the nodule as benign or malignant. The performance of both the CNN and classifiers were similar, with the classifiers performing slightly better. Compared to the performance of radiologists, the sensitivity of nodule detection was within the range

of radiologists at 65% with the two stage neural networks vs 51-81.3% with radiologists. The false positive rate is much higher than the neural networks which is at 6.78 false positives per case with the neural networks vs 0.33-1.39 false positives per case with radiologists. Despite the large number of false positives rate, by solely using the largest nodule detected for cancer prediction. The precision with the classifiers is substantially higher at 41% compared to 1-2% by radiologists.

REFERENCES

- [1] Smita Raut¹, Shraddha Patil², Gopichand Shelke², Lung Cancer Detection using Machine Learning Approach”, International Journal of Advance Scientific Research and Engineering Trends(IJASRET),2021.
- [2] N.Camarlinghi, “Automatic detection of lung nodules in computed tomography images: Training and validation of algorithms using public research databases”, Eur. Phys. J. Plus, vol. 128, no. 9, p. 110, Sep. 2013.
- [3] R. L. Siegel, K. D. Miller, and A. Jemal, “Cancer statistics, 2016”, CA, Cancer J. Clin., vol. 66, no. 1, pp. 730, 2016.
- [4] Detecting and classifying nodules in Lung CT scans, <http://modelheelephant.blogspot.com/2017/11/detecting-and-classifying-nodulesin.html>,2017. [5] Diego Riquelme and Moulay A. Akhloufi, “Deep Learning for Lung Cancer Nodules Detection and Classification in CT Scans”,www.mdpi.com,2020.
- [6] Anita Chaudhary, Sonit Sukhraj Singh, “Lung Cancer Detection on CT Images by using Image Processing”,IEEE,2012.
- [7] Gawade Prathamesh Pratap, R.P. Chauhan, “Detection of Lung Cancer Cells using Image Processing Techniques”, International Conference on Power Electronics, Intelligent Control and Energy Systems(ICPEICES),2016.

- [8] Pooja R. Katre, Dr. Anuradha Thakare, "Detection of Lung Cancer Stages using Image Processing and Data Classification Techniques", International Conference for Convergence in Technology, IEEE, 2017
- [9] Rituparna Sarma, Yogesh Kumar Gupta "A comparative study of new and existing segmentation techniques", ICCRDA, 2020.
- [10] Eali Stephen Neal Joshua^{1*}, Midhun Chakkravarthy¹, Debnath Bhattacharyya², "An Extensive Review on Lung Cancer Detection Using Machine Learning Techniques: A Systematic Study", International Information and Engineering Technology Association (IIETA), 2020.
- [11] Pankaj Nanglia, Sumit Kumar, Aparna N. Mahajan, Paramjit Singh, Davinder Rathee, "A hybrid algorithm for lung cancer classification using SVM and Neural Networks", The Korean Institute of Communication and Information Science (KICS), 2020. Also available at www.elsevier.com/locate/ictc.
- [12] Chao Zhang, Xing Sun, Kang Dang et al "Toward an Expert Level of Lung Cancer Detection and Classification Using a Deep Convolutional Neural Network", The Oncologist, 2019. Also available at www.TheOncologist.com.
- [13] Nidhi S. Nadkarni and Prof. Sangam Borkar, "Detection of Lung Cancer in CT Images using Image Processing", Proceeding of the Third International Conference on Trends and Informatics (ICOEI), IEEE, 2019.
- [14] Ruchita Tekade, Prof. DR. K. Rajeswari, "Lung Cancer Detection and Classification using Deep Learning", Fourth International Conference on Computing Communication Control and Automation (ICCUBEA), IEEE, 2018
- [15] Moffy Vas, Amita Dessai, "Lung Cancer detection system using lung CT image processing", IEEE, 2017