

A Survey on Lung Cancer Diagnosis Research

Minaj Chaugule, Kharat Govind U

Abstract— Lung cancer seems to be the common cause of death among people throughout the world. Early detection of lung cancer can increase the chance of survival among people. The overall 5-year survival rate for lung cancer patients increases from 14 to 49% if the disease is detected in time. There are various research done on lung cancer diagnosis. This paper will review published work and mathematical parameters used in lung cancer diagnosis. This paper will be helpful to students, academicians and researchers to understand scope of further research.

Index Terms—Lung cancer, Mean square error, peak signal to noise ratio, Variance.

I. INTRODUCTION

With respect to engineering field, it becomes challenging to do research over lung cancer diagnosis. But still there are number of papers available, researching on lung cancer diagnosis. Especially using image processing. According to a statistics conducted by world health organization that every year more than 7.6 million people died of lung cancer. Moreover, the death rates of lung cancer are expected upon to keep rising, to wind up around 17 million worldwide in 2030[6]. We found that lung cancers deaths in Bangladesh reached 9,660 or 1.33% of total deaths, according to the latest WHO data published. In year of 2005, around 1,362,825 new cancer cases are expected and around 571,590 deaths are expected to happen due to cancer in the United States. It was evaluated that there will be 162,921 deaths from lung cancer, which occurs 30% of all cancer deaths [7]. Hence Authors decided to do survey on Published work relevant to Lung cancer diagnosis. In next section we will review Papers published on the same topic.

II. PAPER SURVEY

In paper Multi-Stage Lung Cancer Detection and Prediction Using Multi-class SVM Classifier “ published in 2018 International Conference on Computer, Communication, Chemical, Material and Electronic Engineering (IC4ME2), authors says that “Recognition and prediction of lung cancer in the earliest reference point stage can be very useful to improve the survival rate of patients. But diagnosis of cancer is one the major challenging task for radiologist. For detecting, predicting and diagnosing lung cancer, an intelligent computer-aided diagnosis system can be very much useful for radiologist. This paper proposed an efficient lung cancer detection and prediction algorithm using multi-class SVM (Support Vector Machine) classifier. Multi-stage classification was used for the detection of cancer. This system can also predict the probability of lung cancer. In every stage of classification image enhancement and

segmentation have been done separately. Image scaling, color space transformation and contrast enhancement have been used for image enhancement. Threshold and marker-controlled watershed based segmentation has been used for segmentation. For classification purpose, SVM binary classifier was used. Our proposed technique shows higher degree of accuracy in lung cancer detection and prediction”.

[1]

In paper ‘Texture Analysis Based Feature Extraction and Classification of Lung Cancer’, published in 2019 IEEE International Conference on Electrical, Computer and Communication Technologies (ICECCT). Authors states that “Lung cancer is most life-threatening disease, treatment of which must be the primary goal throughout scientific research. The early recognition of cancer can be helpful in curing disease entirely. There are numerous techniques found in literature for detection of lung cancer. Several investigators have contributed their facts for cancer prediction. These papers largely pact about prevailing lung cancer detection techniques that are obtainable in the literature. A numeral of methodologies has been originated in cancer detection methodologies to progress the efficiency of their detection. Diverse applications like as support vector machines, neural networks, image processing techniques are extensively used in for cancer detection which is elaborated in this work.”[2]

In paper ‘Semi-Supervised Multi-Task Learning for Lung Cancer Diagnosis’ published in IEEE 2018, it states that “Early detection of lung nodules is of great importance in lung cancer screening. Existing research recognizes the critical role played by CAD systems in early detection and diagnosis of lung nodules. However, many CAD systems, which are used as cancer detection tools, produce a lot of false positives (FP) and require a further FP reduction step. Furthermore, guidelines for early diagnosis and treatment of lung cancer are consist of different shape and volume measurements of abnormalities. Segmentation is at the heart of our understanding of nodules morphology making it a major area of interest within the field of computer aided diagnosis systems. This study set out to test the hypothesis that joint learning of false positive (FP) nodule reduction and nodule segmentation can improve the computer aided diagnosis (CAD) systems’ performance on both tasks. To support this hypothesis we propose a 3D deep multi-task CNN to tackle these two problems jointly. We tested our system on LUNA16 dataset and achieved an average dice similarity coefficient (DSC) of 91% as segmentation accuracy and a score of nearly 92% for FP reduction. As a proof of our hypothesis, we showed improvements of segmentation and FP reduction tasks over two baselines. Our results support that joint training of these two tasks through a multi-task learning approach improves system performance on both. We also showed that a semi-supervised approach can

be used to overcome the limitation of lack of labeled data for the 3D segmentation task.”[3]

In paper ‘Sex and Smoking Status Effects on the Early Detection of Early Lung Cancer in High Risk Smokers using an Electronic Nose’, published in IEEE Transactions on Biomedical Engineering, 2015. It states that “Objective: Volatile Organic Compounds (VOC) in exhaled breath as measured by electronic nose (e-nose) have utility as biomarkers to detect subjects at risk of having lung cancer in a screening setting. We hypothesize that breath analysis using an e-nose chemo-resistive sensor array could be used as a screening tool to discriminate patients diagnosed with lung cancer from high-risk smokers.

Methods: Breath samples from 191 subjects – 25 lung cancer patients and 166 high-risk smoker control subjects without cancer – were analyzed. For clinical relevancy, subjects in both groups were matched for age, sex, and smoking histories. Classification and Regression Trees and Discriminant Functions classifiers were used to recognize VOC patterns in e-nose data. Cross-validated results were used to assess classification accuracy. Repeatability and reproducibility of e-nose data were assessed by measuring subject-exhaled breath in parallel across two e-nose devices. Results: E-nose measurements could distinguish lung cancer patients from high-risk control subjects, with a better than 80% classification accuracy. Subject sex and smoking status impacted classification as area under the curve results (ex-smoker males 0.846, ex-smoker female 0.816, current smoker male 0.745 and current smoker female 0.725) demonstrated. Two e-nose systems could be calibrated to give equivalent readings across subject-exhaled breath measured in parallel. Conclusions: E-nose technology may have significant utility as a non-invasive screening tool for detecting individuals at increased risk for lung cancer.

Significance: The results presented further the case that VOC patterns could have real clinical utility to screen for lung cancer in the important growing ex-smoker population.”[4]

In paper ‘Study of Malignancy Associated Changes in Sputum Images as an Indicator of Lung Cancer’, published in Proceedings of the 2016 IEEE Students’ Technology Symposium .it states that “Lung cancer is one among the major causes of cancer related deaths. Fortunately, an early stage diagnosis can increase the survival rates of the patients. Sputum cytology is one of the easiest and cost-effective method for lung cancer diagnosis. Chances of misdiagnosis and sampling error related to sputum cytology led to the concept of malignancy associated changes. Malignancy associated changes (MAC) are the subtle changes that happens to the normal appearing cells near or distant from the malignant cells. Literature suggests that these changes can be used as an indicator for lung cancer rather than using malignant cells which are very less in number compared to the normal appearing cells in sputum cytology images. The proposed work is intended to detect cells with MAC from sputum smear images. Analysis of nuclei texture features of sputum cell nuclei using Gray Level Co-occurrence Matrix and Gray Level Run Length Matrix from both normal and cancer patients revealed that both type of cells could be differentiated. Among 110 texture features calculated for each nuclei, a set of 35 features which clearly distinguishes normal cells and normal appearing cells were chosen. Support Vector Machine (SVM) classifier is used to classify the cells into two classes’ i.e cells with MAC and cells without MAC.

This study demonstrates that the presence of MAC cells in conventional microscopic sputum cytology Images can be identified using image processing techniques and it can have some significance in the early detection of lung cancer.”[5]

In next section we will review some mathematical parameters used in lung cancer diagnosis.

III. MATHEMATICAL PARAMETERS

Mathematical parameters are key of lung cancer diagnosis. Based on mathematical parameters only diagnosis success is dependent. In this section we will see those important mathematical parameters.

2.1 Statistical Analysis

a) Entropy: It indicates average information of the image. The lowest value of entropy means no uncertainty of the image information. It is zero if the event is sure or impossible

$$E = -\sum_x \sum_y P[x, y] \log P[x, y] \quad (1)$$

$P[x, y]$ Is the probability difference between two adjacent pixels and log is the base2 logarithm. Deliberating Entropy $E=0$ if $P=0$ or 1. Entropy is supposed to be high throughout the image and is calculated by equation (1).

b) Mean: It calculates the mean of the gray levels in the image. The Mean is supposed to be high, if the sum of the gray levels of the image is high. Mean depends on the first moment of the data. Technically, a moment is defined by a mathematical formula that just so happens to equal formulas for some measures in statistics. First moment is the mean which is represented as in equation (2) and (3)

$$S^{th} = \frac{(x_1^s + x_2^s + x_3^s + \dots + x_n^s)}{n} \quad (2)$$

First moment (S=1)

$$S^{th} = \frac{(x_1^1 + x_2^1 + x_3^1 + \dots + x_n^1)}{n} \quad (3)$$

This formula is identical to the formula to find the sample mean. Just add up all of the values and divide by the number of items in given data set. The mathematical expression of mean is given as in equation (4)

$$\mu = 1/N * M \sum_{x=0}^M \sum_{y=0}^N P[x, y] \quad (4)$$

Where $N * M$ (255*255) is the size of the image

c) Variance: It explains about the distribution of gray levels over the image. The value of the Variance is expected to be high, if the gray levels of the image are spread out extensively. It explains about the probability of distribution, describing how far the value lies from the mean that is anticipated value, which can also be defined as the moments of a distribution. Second moment (S=2) is given as by equation (5). The second moment is the Variance. μ_x is the average of x . Mean provides each pixel intensity for the whole image, whereas the variance gives each pixel variations from the neighboring pixels and is use to classify image into different regions or areas. It is the average of the squared differences from the mean. It is the variability around the value.

$$S^{th} = \sum (x_i - \mu_x)^2 \quad (5)$$

Steps to calculate variance of an image

1. Calculate Mean (simply average of numbers)
2. For each number, subtract the mean and the square of result (the squared difference)
3. Average of those squared differences

The mathematical expression for calculating Variance is given in equation (6), $N-1$ can be changed to N if the \bar{x} is known prior rather than being estimated from the data

$$\text{var} = \sigma^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (6)$$

d) Standard Deviation: If the value of standard deviation is less which means that the majority of the data is near the mean value and if it is more that means data is more spreaded over the image. The value of standard deviation is assigned to the center pixel of the image. All the steps are similar to calculate the standard deviation as the variance, only the last step is added as the square root, hence is the square root of the variance represented by equation (7)

$$SD = \sigma = \sqrt{\sigma^2} \quad (7)$$

e) MSE (Mean Square Error) : The MSE represents the averaging of the squares of the errors between the two images [9]. The error is the amount by which the values of the reference image differ from test image calculated form equation (8).

$$MSE = \sum_0^{m-1} \sum_0^{n-1} \|f(i, j) - g(i, j)\|^2 \quad (8)$$

$f(i, j)$ Represents the matrix data of original image and $g(i, j)$ represents the matrix data of test image. m represents the numbers of rows of pixels of the images and i represents the index of that row n represents the number of columns of pixels of the image and j represents the index of that column. MSE for the practical purpose allows comparing the true pixel values of original image to cancerous image.

f) Correlation: Correlation is an additionally a statistical procedure which demonstrates how factors are robustly related with each other. It extracts necessary information from the image. It is used to find the location in an image that is analogous to the reference image. Correlation is a measure of gray level linear dependence between the pixels at the specified positions relatively [10].

$$y(n) = \sum_{-n}^n x(v)h(n-v) \quad (9)$$

Where $x(v)$ - Image1 and $h(n-v)$ -Image 2(Shifted).

IV. CONCLUSION

A survey on lung cancer diagnosis has large scope. But in this paper we tried to do survey on one of its significant elements, like published work and mathematical parameters. Published work is essential one as it gives us idea about pros and cons of existing work. It helps researchers, students and

academicians to know further scope for work. As well as mathematical parameters is also very important element of lung cancer diagnosis. Based on these parameters only algorithms and techniques are built for successful diagnosis. Here we covered both elements. In future we can do survey on classification and segmentation techniques.

REFERENCES

- [1] Ojala, T., Pietikainen, M., & Maenpaa, T. (2002). Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *IEEE Transactions on pattern analysis and machine intelligence*, 24(7), 971-987.
- [2] Radüntz, T., Scouten, J., Hochmuth, O., & Meffert, B. (2015). EEG artifact elimination by extraction of ICA-component features using image processing algorithms. *Journal of neuroscience methods*, 243, 84-93.
- [3] Sairamya, N. J., Selvaraj, T. G., Ramasamy, B., Deivendran, N. P., & Subathra, M. S. P. (2018). Classification of EEG signals for detection of epileptic seizure activities based on feature extraction from brain maps using image processing algorithms. *IET Image Processing*.
- [4] Sairamya, N. J., George, S. T., Subathra, M. S. P., & Kumar, N. M. (2019). Computer-Aided Diagnosis of Epilepsy Based on the Time-Frequency Texture Descriptors of EEG Signals Using Wavelet Packet Decomposition and Artificial Neural Network. In *Cognitive Informatics and Soft Computing* (pp. 677-688). Springer, Singapore.
- [5] Sairamya, N. J., George, S. T., Ponraj, D. N., & Subathra, M. S. P. (2017, October). Automated Detection of Epileptic Seizure Using Histogram of Oriented Gradients for Analysing Time Frequency Images of EEG Signals. In *International Conference on Next Generation Computing Technologies* (pp. 932-943). Springer, Singapore.
- [6] Kaya, Y., Uyar, M., Tekin, R., & Yıldırım, S. (2014). 1D-local binary pattern based feature extraction for classification of epileptic EEG signals. *Applied Mathematics and Computation*, 243, 209-219.
- [7] Jaiswal, A. K., & Banka, H. (2017). Local pattern transformation based feature extraction techniques for classification of epileptic EEG signals. *Biomedical Signal Processing and Control*, 34, 81-92.
- [8] Tiwari, A. K., Pachori, R. B., Kanhangad, V., & Panigrahi, B. K. (2017). Automated diagnosis of epilepsy using key-point based local binary pattern of EEG signals. *IEEE journal of biomedical and health informatics*, 21(4), 888-896.
- [9] Sairamya, N. J., George, S. T., Ponraj, D. N., & Subathra, M. S. P. (2018). Detection of Epileptic EEG Signal Using Improved Local Pattern Transformation Methods. *Circuits, Systems, and Signal Processing*, 1-22.
- [10] Sairamya, N. J., George, S. T., Balakrishnan, R., & Subathra, M. S. P. (2018). An effective approach to classify epileptic EEG signal using local neighbor gradient pattern transformation methods. *Australasian physical & engineering sciences in medicine*, 41(4), 1029-1046.

First Author Minaj Chaugule has done her Bachelor of Engineering in Electronica and telecommunication. Now she is pursuing her Masters of Engineering in Embedded and VLSI.
 Email: minajchaugule08@gmail.com

Second Author Prof. Kharat Govind U has done his Bachelor of Engineering in Electronics and Masters of Engineering in Power systems. His PhD topic was power Electronics Now he is working as a Principal at Sharadchandra college of engineering, Otur, Maharashtra, India.
 Email: principalspc09@gmail.com