# LUNG CANCER PREDICTION USING GLCM WITH HYPERPARAMETER OPTIMIZATION APPROACH

**Prakasha Raje Urs M[1], Dr. G N K Suresh Babu[2]**
[1]Research Scholar, Srishti College of Commerce and Management, Banashankari, Bengaluru, Affiliated to University of Mysore, India, 560085.
[2]Professor, Department of Computer Science, Srishti College of Commerce and Management, Banashankari, Bengaluru, Affiliated to University of Mysore, India, 560085.

**Abstract:** *Cancer is one of the most lethal illness professionals. This is responsible for an uncountable number of deaths throughout the globe.The concept of diagnosing lung cancer at an earlier stage piqued the interest of medical professionals.This research presents a novel method for the detection of lung cancer that is based on the processing of pictures obtained from CT scans.With the use of cases from the Lung Imaging Database Consortium (LIDC) database, we were able to assess the practicability of applying algorithms for the detection of lung cancer in this research. The primary purpose of this study is to determine if the tumours that are discovered in the lung are malignant or benign. This has been achieved using the GLCM feature extractor and the SVM classifier. The proposed research includes the hyperparameter optimization approach that has been proposed for use in identifying lung cancer in its early stages.* Results obtained with hyperparameter optimization has showed the overall accuracy prediction of 92.06% and more than 90% of accuracy in deciding *tumour as malignant, benign, or normal*. And area under the curve (AUC) value 0.993,0.996 and0.997 has been obtained for the classes*malignant, benign, or normal*respectively during the classification cancer.
Keywords: Area Under Curve, Benign, Hyperparameter, Lung Cancer, Malignant, Optimization.

## INTRODUCTION:

*According to WHO (World Health Organization) lung cancer contributes about 14 per cent among all the cancers. Therefore, early detection and treatment is very much required. Computed Tomography (CT) scan can provide valuable information in the diagnosis of lung diseases.* Lung cancer is a fatal illness that affects a significant number of people worldwide. Eachyear, lung cancer takes more lives than all new types of cancer combined, including brain, prostate, and breast [1]. Lung cancer is one of the most common cancers, accounting for over 225,000 cases, 150,000 deaths, and $12 billion in health care costs yearly in the U.S. [2]. It is also one of the deadliest cancers; overall, only 17% of people in the U.S. diagnosed with lung cancer survive five years after the diagnosis, and the survival rate is lower in developing countries. The stage of a cancer refers to how extensively it has metastasized [3]. Stages 1 and 2 refer to cancers localized to the lungs and latter stages refer to cancers that have spread to other organs. Current diagnostic methods include biopsies and imaging, such as CT scans. Early detection of lung cancer (detection during the earlier stages) significantly improves the chances for survival, but it is also more difficult to detect early stages of lung cancer as there

are fewer symptoms [4]. Lung cancer kills more people every year than colon, breast, andprostate cancer integrated, accounting for about 25% of cancer-related fatalities. Numerous sophisticated methods, including chest radiography (X-ray), computed tomography (CT), sputumcytology, and magnetic resonance imaging (MRI), are used to diagnose lung cancer in its earliest stages [5-6]. Therefore, another advancement is crucial for the early identification of lungcancer. Individuals over the age of 40 are disproportionately affected [7-8]. According to theWorld Health Organization, lung cancer kills over 7.6 million lives each year. In comparison to various types of cancer, lung cancer takes the most significant number oflives. Lung cancer accounts for around one-quarter of all cancer-related fatalities. The two kindsof lung cancer are non-small cell and small-cell lung cancer [9]. It consists of four stages (Fig.1). Cancer begins in the lungs during the first stage and spreads to the chest during the secondaryand third stages [10-11]. By stage 4, it has spread throughout the body. Individuals who smoke asmuch as possible enhance their chance of transmitting the infection. Different procedures, suchas the X-bar, engaged resonation imaging scan, and computed tomography scan, are utilised infacilities to distinguish cancer at an early stage [12-13]. The fundamental downside of these solutions isthe higher expense and work involved in disclosing them. Additionally, doctors generally assesscancer patients based on their symptoms, which commonly manifest during the disease's latest ages [14-15]. Lung cancer has a long-term survival rate of 15%. Therefore, early detectionof lung cancer is crucial for increasing survival chances. Convolutional neural networks have achieved better than Deep Belief Networks in current studies on benchmark computer vision datasets [16]. The CNNs have attracted considerable interest in machine learning since they have strong representation ability in learning useful features from input data in recent years. In this paper, we apply an extensive pre-processing technique to get the accurate nodules to enhance the accuracy of detection of lung cancer [17-18]. Moreover, we perform an end-to-end training of CNN from scratch to realize the full potential of the neural network i.e., to learn discriminative features. Extensive experimental evaluations are performed on a dataset comprising lung nodules from more than 1390 low dose CT scans. The diagnosis of lung cancer based on the microscopic examination of lung tissue samples has several limitations [19-20].One of these is that doctors still rely on their own individual interpretations of what they see.When it comes to diagnosing lung cancer in patients, a medical practitioner must do careful observation and precise analysis.Thus, there is a need for a system that is capable of autonomously diagnosing lung cancer from microscopic pictures of biopsy samples. Using digital image processing and analysis of microscopic images obtained from biopsies, the purpose of this study is to develop a method for the early diagnosis of lung cancer.

## METHODOLOGY

Typical CAD systems for lung cancer have the following pipeline: image pre-processing, detection of cancerous nodule, nodule false positive reduction, malignancy prediction for each nodule, and malignancy prediction for overall CT scan [15]. The proposed methodology has four phases for the classification of lung cancer. In phase one, the required data is collected from the database https://www.kaggle.com/datasets. In phase two resize of image has been done. In phase three features are extracted using GLCM (Gray Level Co-occurrence Matrix). These extracted features are used in phase four for classification purpose which is carried by
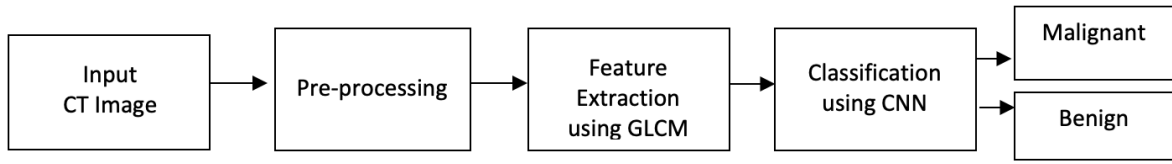
CNN (Convolutional Neural Network)



**Figure 1: Block diagram of Lung cancer prediction**

This paper presents lung cancer detection based on chest CT images using CNN. The first stage involves extracting lung regions from the CT picture and segmenting each slice in those regions to find malignancies. The CNN architecture is trained using the segmented tumour regions. The dataset is divided into three sets for the training, validation, and testing phases in the ratio 70%:20%:10% after the images are ready in their binary matrix format. The patient images are then evaluated using CNN.

**Feature Extraction**

Feature extraction is applied to extract features like energy, entropy, variance. Normality and abnormality of an image can be determined in this stage. Various features of an image can be Contrast, Energy, Homogeneity, Mean, Standard Deviation, Entropy, RMS, Variance, Smoothness etc. Each object is extracted for its features based on certain parameters and is then assigned a certain class. The GLCM feature extraction method is a matrix that describes the occurrence frequency of two pixels. Following equations are used to extract the features out of the lung images.

| 1 | Contrast | $\sum_i \sum_j (i-j)^2 p_d(i,j)$ |
|---|----------|---------------------------------|
| 2 | Energy | $Energy = \sqrt{ASM}$<br>$ASM = \sum_i \sum_j p_d{}^2(i,j)$ |
| 3 | Entropy | $-\sum_i \sum_j p_d(i,j)\, lnln\, p_d(i,j)$ |
| 4 | Homogeneity | $\sum_i \sum_j \dfrac{1}{1+(i-j)^2} p_d(i,j)$ |
| 5 | Sum of entropy | $-\sum_{i=2}^{2N_g} p_{x+y}(i) \log\log \{p_{x+y}(i)\} = f_8$ |
| 6 | Sum of variance | $\sum_{i=2}^{2N_g} (i-f_8)^2\, p_{x+y}(i)$ |
| 7 | Dissimilarity | $\sum_{j=1}^{N} |i-j| \cdot p(i,j)$ |
| 8 | Sum of average | $\sum_{i=2}^{2N_g} i p_{x+y}(i)$ |

**SVM classifier**

Although Support Vector Machines (SVM) is typically thought of as a classification tool, it can also be used to solve regression problems. Many continuous and categorical variables can

be handled with ease. To divide various classes, SVM creates a hyper plane in multidimensional space. SVM generates ideal hyper plane in an iterative manner, which is utilised to minimise an error. Finding a maximum marginal hyper plane (MMH) that optimally classifies the dataset is the central goal of SVM. SVM's primary benefit is that it can be applied to both classification and regression issues. Here, we have developed SVM classifier for multiclass application. In the proposed method SVM must classify the lung CT images into benign, malignant and Normal class.
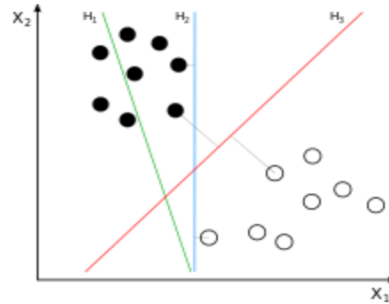


**Figure 2:  Block diagram of Lung cancer prediction**

The above-mentioned approach ensures that with larger margins, the classifier yields a lower generalization error.  Separation hyper planes as shown in Figure 2, in which H1 does not divide the two classes; H2 divides them but with a small separating margin while H3 separates the classes with the best margin.

**RESULTS**

For the computation processes we have considered the Lung Image Database Consortium and Image Database Resource Initiative (LIDC-IDRI) database, provide training data set. The size of CT image will be 512x512x3. At the time of feature extraction this image of both the sets are resized to 224x224x3. Here, 500 images have been used to conduct the experiment. Out of the available data, 70% of image data for training the model and other 30% for verifying the result and for checking accuracy of the network. Our research has carried out with and without hyper parameter tuning process to predict the possibilities of the disease in the lung images. Statistical parameters like True Positive Rate (TPR), False Positive Rate, Accuracy and Confusion matrix have been computed and tabulated to measure the performance of the system. Here, class 1 indicates benign, class 2 indicates malignant and class 3 as normal case condition.

**Without Hyperparameter Tuning:**

The proposed method has carried out using the combination of GLCM and SVM with and without hyper parameter tuning.  The extracted features from GLCM are serve as inputs for the SVM classifier, whereas the outputs are categories of Benign (1), Malignant (2) and Normal (3). In Table 1 prediction of these classes for first 10 cases has been tabulated and it shows that in few cases class 2 which malignant is wrongly predicted as class 1 which is Benign and class 3 which is Normal.

**Table 1: Prediction chart Without Hyperparameter Tuning**

```
TrueLabel      PredLabel              Posterior

   2            {'1'}      0.50919     0.32321      0.1676
   2            {'2'}      0.10901     0.87188      0.019109
   2            {'1'}      0.50919     0.32321      0.1676
   2            {'2'}      0.11797     0.87249      0.0095412
   2            {'2'}      0.11885     0.8716       0.0095529
   3            {'3'}       0.4124     0.1327       0.4549
   2            {'3'}      0.31966     0.22883      0.45152
   2            {'2'}       0.0151     0.96814      0.016756
   2            {'2'}      0.21438     0.68082      0.10479
   2            {'2'}      0.20042     0.68341      0.11617
```

The proposed model overall accuracy has been calculated and tabulated here considering the 3 types like benign, malignant, and normal. In developing any system, k-fold cross-validation helps us to build the model as a generalized one. Here, 5-fold cross validation has been achieved and results corresponding to each fold is evaluated and tabulated as in Table 2 and Table 3.

**Table 2: Statistical Parameter analysis Without Hyperparameter Tuning**

| Fold | Case | FPR | TPR | Precision | Recall | F1-Score | Accuracy |
|------|------|-----|-----|-----------|--------|----------|----------|
| 1 | B | 0.018 | 0.981 | 0.973 | 0.939 | 0.955 | 0.963 |
| | M | 0.027 | 0.973 | 0.942 | 0.931 | 0.936 | 0.960 |
| | N | 0.040 | 0.960 | 0.895 | 0.958 | 0.987 | 0.960 |
| 2 | B | 0.051 | 0.950 | 0.929 | 0.912 | 0.920 | 0.934 |
| | M | 0.025 | 0.9752 | 0.925 | 0.886 | 0.905 | 0.952 |
| | N | 0.033 | 0.968 | 0.935 | 0.989 | 0.944 | 0.974 |
| 3 | B | 0.042 | 0.958 | 0.937 | 0.981 | 0.958 | 0.967 |
| | M | 0 | 1 | 1 | 0.927 | 0.962 | 0.978 |
| | N | 0.016 | 0.984 | 0.965 | 0.976 | 0.956 | 0.982 |
| 4 | B | 0.047 | 0.953 | 0.942 | 0.912 | 0.926 | 0.934 |
| | M | 0.005 | 0.995 | 0.984 | 0.913 | 0.947 | 0.975 |
| | N | 0.062 | 0.938 | 0.863 | 0.962 | 0.965 | 0.945 |
| 5 | B | 0.096 | 0.904 | 0.871 | 0.878 | 0.874 | 0.894 |
| | M | 0.025 | 0.975 | 0.926 | 0.840 | 0.881 | 0.938 |
| | N | 0.047 | 0.953 | 0.898 | 0.963 | 0.994 | 0.956 |

Here, figure 3 and figure 4 represents result pertaining to one of the fold out of 5-fold validation outcome. Confusion matrix evaluates the performance of the classification models, when they make predictions on test data, and tells how good our classification model is. It not only tells the error made by the classifiers but also the type of errors such as it is either benign (B), malignant (M) and Normal (N). Here, our model is providing an highest accuracy of 94.7% in detecting the levels like Normal and with a lowest accuracy of 86.6% while detecting Malignant. overall accuracy of the system will results in 93.1% towards detection or prediction DR.
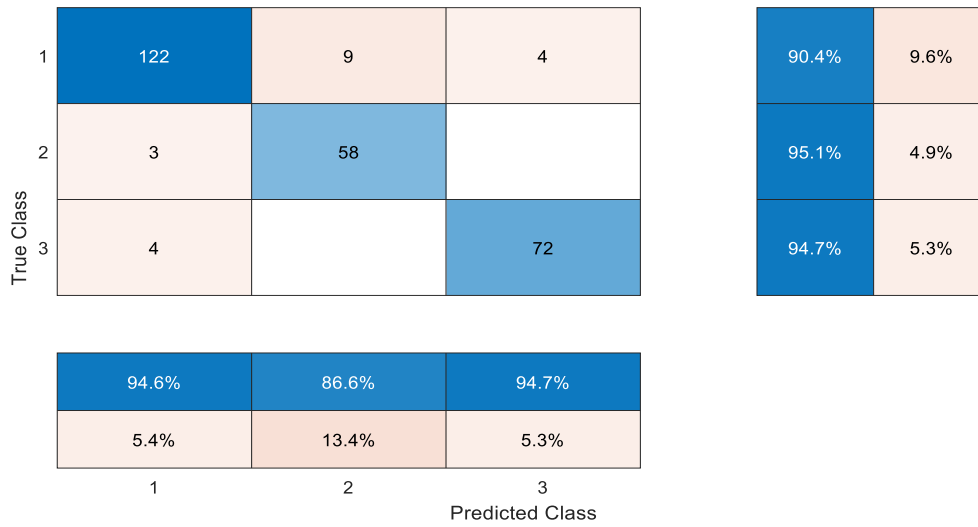
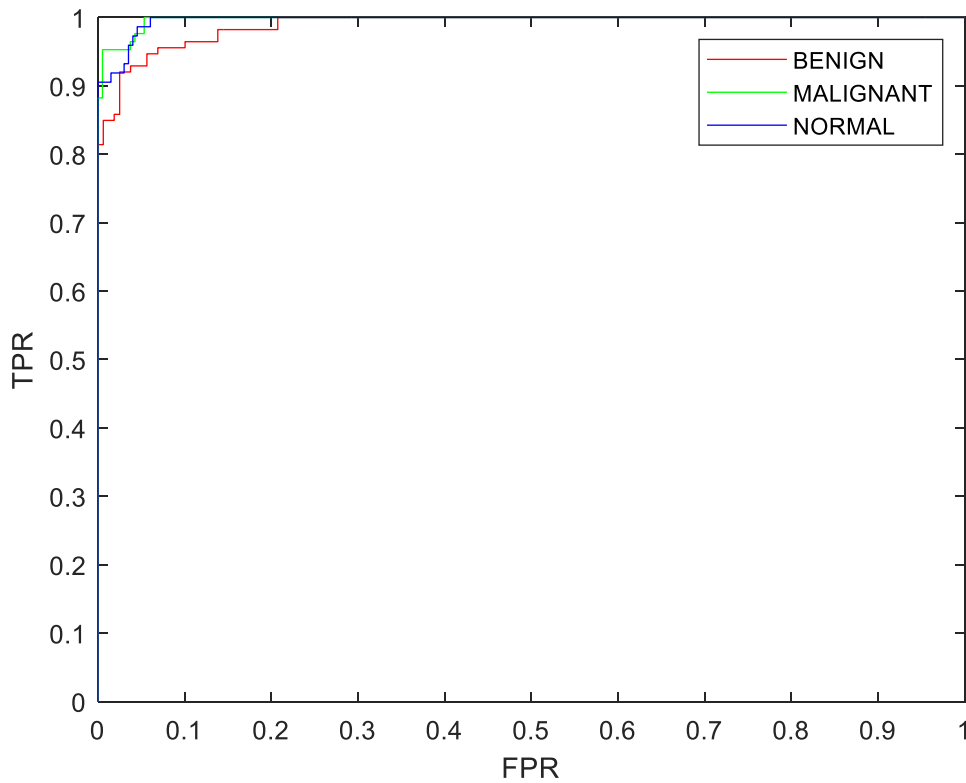**Figure 3: Confusion Matrix Without Hyperparameter Tuning**



**Figure 4: ROC plot for Without Hyperparameter Tuning**

The ROC curve and the corresponding AUC for each class is as given in Figure 4 and Table 2 respectively. 5- Fold cross validation results have been recorded in Table 3 to exhibit AUC value for the each of the class. From the table it concludes that maximum AUC value 0.998 has obtained for class malignant as well as for normal.

**Table 3: AUC for each class Without Hyperparameter Tuning**

| Fold | Benign | Malignant | Normal |
|------|--------|-----------|--------|
| 1 | 0.989 | 0.997 | 0.996 |
| 2 | 0.993 | 0.998 | 0.998 |
| 3 | 0.997 | 0.998 | 0.997 |
| 4 | 0.986 | 0.995 | 0.997 |
| 5 | 0.995 | 0.995 | 0.998 |

**With Hyperparameter Tuning:**

Hyperparameter minimizes the five-fold cross-validation loss by using automatic hyperparameter optimization. Here, max Objective Evaluations of 30 reached. The optimization minimizes the cross-validation loss (error) using by varying the parameters. A typical loss curve corresponding to hyper parameter tuning is as in figure 5.
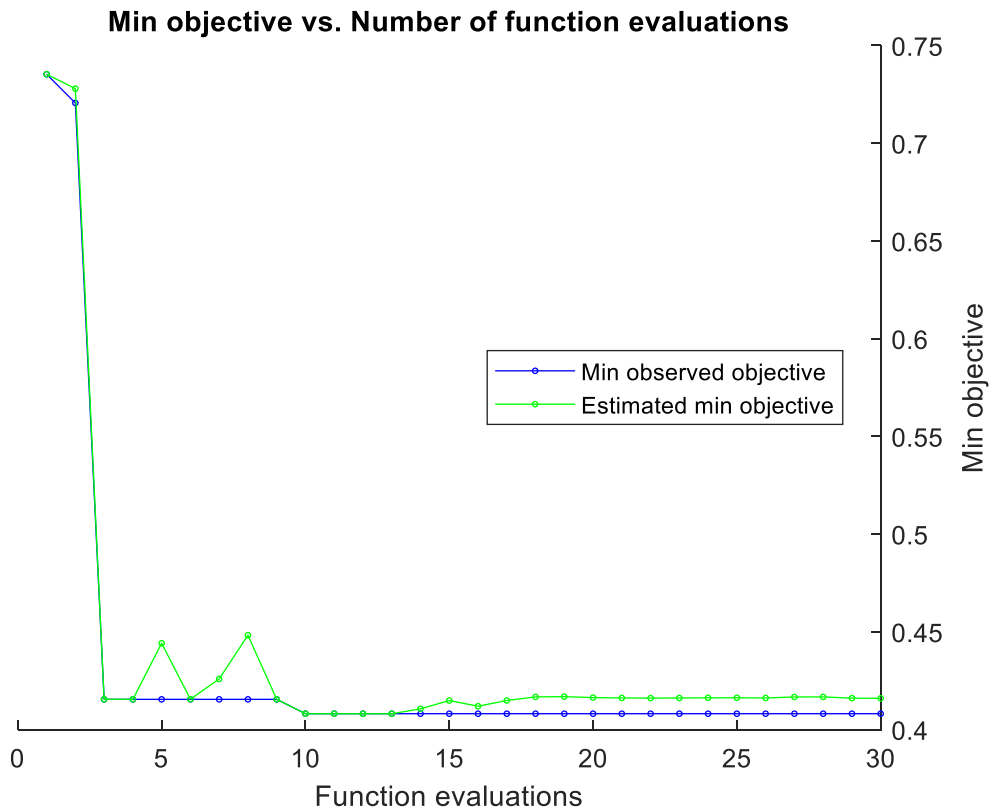


**Figure 5: Loss curve for With Hyperparameter Tuning**

Table 4 indicates chart of 30 evaluation objectives for which it shown best or acceptable as evaluation result. Predictions from all learners have been computed. Loss for all observations has been computed and recorded as in figure 5.

**Table 4: Evaluation of Objectives Without Hyperparameter Tuning**

| Iter | Eval result | Objective | Objective runtime | BestSoFar (observed) | BestSoFar (estim.) |
|---|---|---|---|---|---|
| 1 | Best | 0.73529 | 24.142 | 0.73529 | 0.73529 |
| 2 | Best | 0.72059 | 15.729 | 0.72059 | 0.72794 |
| 3 | Best | 0.41544 | 7.3502 | 0.41544 | 0.41547 |
| 4 | Accept | 0.72426 | 14.001 | 0.41544 | 0.41547 |
| 5 | Accept | 0.53309 | 0.14149 | 0.41544 | 0.44414 |
| 6 | Accept | 0.45956 | 8.8648 | 0.41544 | 0.41552 |
| 7 | Accept | 0.4375 | 9.1631 | 0.41544 | 0.42585 |
| 8 | Accept | 0.47426 | 2.5451 | 0.41544 | 0.44824 |
| 9 | Accept | 0.62868 | 9.2165 | 0.41544 | 0.41547 |
| 10 | Best | 0.40809 | 6.7782 | 0.40809 | 0.40806 |
| 11 | Accept | 0.40809 | 8.5087 | 0.40809 | 0.40807 |
| 12 | Accept | 0.41544 | 6.09 | 0.40809 | 0.40805 |
| 13 | Accept | 0.42279 | 7.7339 | 0.40809 | 0.40806 |
| 14 | Accept | 0.41544 | 8.157 | 0.40809 | 0.41059 |
| 15 | Accept | 0.43015 | 8.1264 | 0.40809 | 0.41482 |
| 16 | Accept | 0.49632 | 0.68552 | 0.40809 | 0.41192 |
| 17 | Accept | 0.43015 | 8.1158 | 0.40809 | 0.41484 |
| 18 | Accept | 0.43015 | 8.2141 | 0.40809 | 0.41669 |
| 19 | Accept | 0.48897 | 11.699 | 0.40809 | 0.4168 |
| 20 | Accept | 0.53676 | 7.7507 | 0.40809 | 0.41633 |
| 21 | Accept | 0.69485 | 13.818 | 0.40809 | 0.41615 |
| 22 | Accept | 0.60294 | 0.095212 | 0.40809 | 0.41606 |
| 23 | Accept | 0.63235 | 0.094357 | 0.40809 | 0.41613 |
| 24 | Accept | 0.63235 | 0.10026 | 0.40809 | 0.41622 |
| 25 | Accept | 0.63235 | 11.184 | 0.40809 | 0.41624 |
| 26 | Accept | 0.70956 | 11.105 | 0.40809 | 0.41614 |
| 27 | Accept | 0.48529 | 11.75 | 0.40809 | 0.41665 |
| 28 | Accept | 0.63235 | 0.10209 | 0.40809 | 0.41669 |
| 29 | Accept | 0.54779 | 7.993 | 0.40809 | 0.41604 |
| 30 | Accept | 0.69853 | 18.619 | 0.40809 | 0.41596 |

Computing posterior probabilities and predict the training-sample labels and class posterior probabilities are presented as in Table 5. Prediction of these classes for 10 cases has been tabulated and it shows that for only one of the cases of class 2 which malignant is wrongly predicted as class 1 which is Benign and class 3 which is Normal. For remaining all other cases it is predicting correctly.

**Table 5: AUC for each class Without Hyperparameter Tuning**

```
TrueLabel      PredLabel                    Posterior
_____      _____      _____

    1            {'1'}           0.53565      0.11699       0.34736
    1            {'1'}           0.65081      0.17032       0.17886
    1            {'1'}           0.54663       0.3052       0.14816
    3            {'3'}           0.25524      0.15997       0.58479
    2            {'2'}         0.00080288     0.99901     0.00018583
    2            {'2'}           0.46499      0.28175       0.25327
    3            {'3'}           0.32718      0.11621        0.5566
    2            {'1'}           0.61381      0.29768      0.088507
    3            {'3'}           0.50549      0.40929       0.08522
    3            {'3'}           0.48022      0.22754       0.29224
```

Using hyperparameter, overall accuracy of the proposed model has been calculated and tabulated as in table 6. Since hyperparameter auto tunes itself during prediction, k-fold cross-validation is not essential in this approach. Using this approach accuracy towards determining all the 3 cases found to be more than 90%. Using these statistical parameters from table 6 confusion matrix and ROC curve has been plotted as in figure 6 and figure 7 respectively. The results out of the plot indicates, the performance of the model is moderately good.

**Table 6: Statistical Parameter for each class With Hyperparameter Tuning**

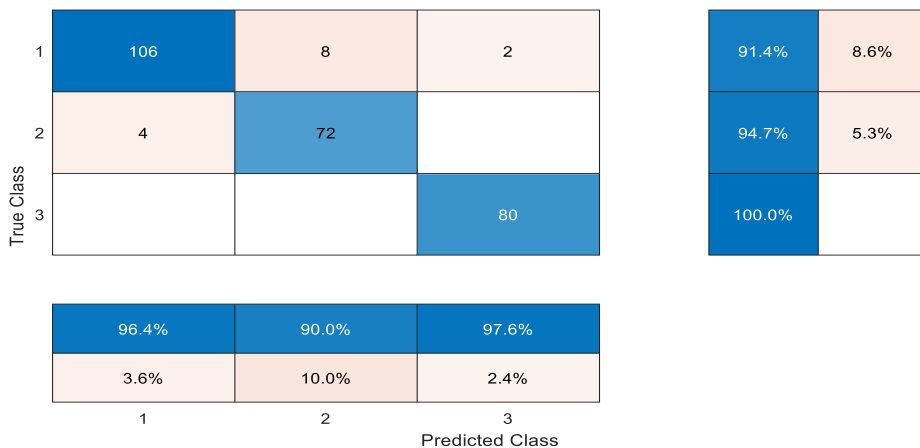| Class | FPR | TPR | Precision | Recall | F1-Score | AUC | Accuracy |
|-------|-----|-----|-----------|--------|----------|-----|----------|
| **Benign** | 0.061 | 0.939 | 0.912 | 0.954 | 0.932 | 0.993 | 0.945 |
| **Malignant** | 0.011 | 0.989 | 0.975 | 0.889 | 0.930 | 0.996 | 0.956 |
| **Normal** | 0.025 | 0.975 | 0.935 | 0.973 | 0..993 | 0.997 | 0.975 |



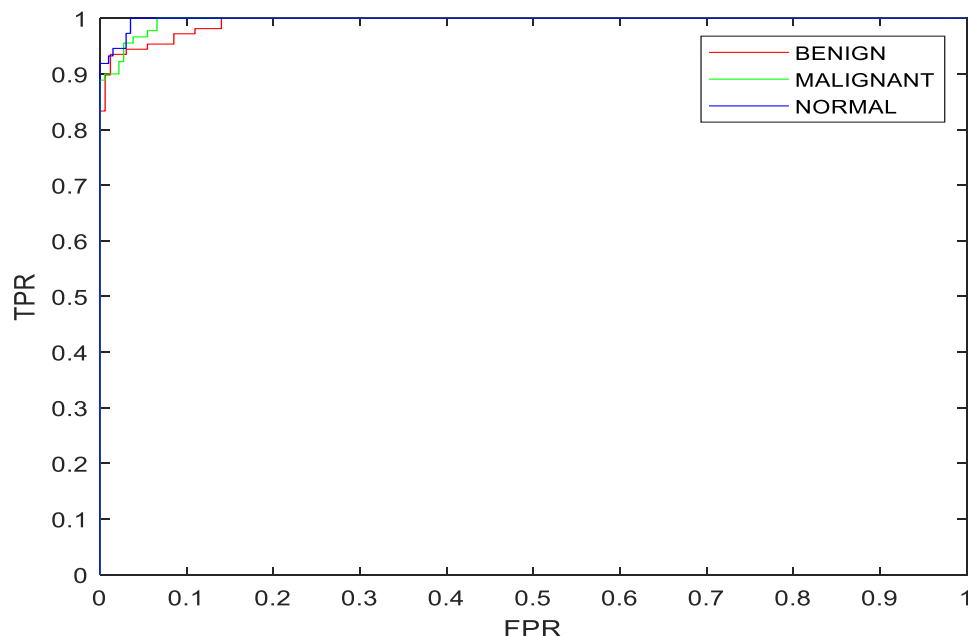**Figure 6: Confusion Matrix for With Hyperparameter Tuning**

**Figure 7: ROC curve for With Hyperparameter Tuning**

## CONCLUSION

To predict lung cancer, this study suggested an efficient CT classification system based on the GLCM and SVM. Experiments to categorise lung cancer into the benign, malignant, and normal classes using standard dataset. When a result, the CT scans tend to vary more considerably amongst patients as the number of patients with lung cancer is found. As technology and computer-aided diagnosis have developed, many automated techniques have been developed to address the problem of lowering false positives when evaluating the presence of cancer in patients' CT images. From the proposed it is found that, we may determine the benign, malignant, and Normal case using hyperparameter tuning with an accuracy of 96.4 %, 90% and 97.6% respectively. Overall accuracy of the system is 94.85%. For the same data set our proposed model With Hyperparameter Tuning performs moderately better than Without Hyperparameter Tuning approach.

## REFERENCES

[1]. He, Jie, Ni Li, W. Q. Chen, Ning Wu, H. B. Shen, Yu Jiang, Jiang Li, Fei Wang, and J. H. Tian. "China guideline for the screening and early detection of lung cancer (2021, Beijing)." Zhonghua Zhong Liu Za Zhi [chinese Journal of Oncology] 43, no. 3 (2021): 243-268.

[2]. Jacobs, Colin, Arnaud AA Setio, Ernst T. Scholten, Paul K. Gerke, Haimasree Bhattacharya, Firdaus A. M. Hoesein, Monique Brink et al. "Deep learning for lung cancer detection on screening CT scans: results of a large-scale public competition and an observer study with 11 radiologists." Radiology: Artificial Intelligence 3, no. 6 (2021): e210027.

[3]. Lee, Jong Hyuk, Dongheon Lee, Michael T. Lu, Vineet K. Raghu, Chang Min Park, Jin Mo Goo, Seung Ho Choi, and Hyungjin Kim. "Deep learning to optimize candidate

selection for lung cancer CT screening: advancing the 2021 USPSTF recommendations." Radiology 305, no. 1 (2022): 209-218.

[4]. Chen, Ke, Lei Liu, Bo Nie, Binchun Lu, Lidan Fu, Zichun He, Wang Li, Xitian Pi, and Hongying Liu. "Recognizing lung cancer and stages using a self-developed electronic nose system." Computers in Biology and Medicine 131 (2021): 104294.

[5]. Chabon, J.J., Hamilton, E.G., Kurtz, D.M., Esfahani, M.S., Moding, E.J., Stehr, H., Schroers-Martin, J., Nabet, B.Y., Chen, B., Chaudhuri, A.A. and Liu, C.L., 2020. Integrating genomic features for non-invasive early lung cancer detection. Nature, 580(7802), pp.245-251.

[6]. Asuntha, A., and Andy Srinivasan. "Deep learning for lung Cancer detection and classification." Multimedia Tools and Applications 79 (2020): 7731-7762.

[7]. Elnakib, Ahmed, Hanan M. Amer, and Fatma EZ Abou-Chadi. "Early lung cancer detection using deep learning optimization." (2020): 82-94.

[8]. Dr. Nagaraju C, Sangana Manasa Durga, Rachana N, "Ascertainment of Lung Cancer at an Early Stage", International Journal of Scientific Research in Computer Science, Engineering and Information Technology, ISSN: 2456-3307,Volume 2, Issue 4, 2017.

[9]. Ignatious S, Joseph R, "Computer Aided Lung Cancer Detection System." IEEE, Proceedings of 2015 Global Conference on Communication Technologies (GCCT 2015), pp. 555–558, 2015.

[10]. Ueda, Daiju, Akira Yamamoto, Akitoshi Shimazaki, Shannon Leigh Walston, Toshimasa Matsumoto, Nobuhiro Izumi, Takuma Tsukioka et al. "Artificial intelligence-supported lung cancer detection by multi-institutional readers with multi-vendor chest radiographs: a retrospective clinical validation study." BMC cancer 21 (2021): 1-8.

[11]. Acheampong, Felix, Trevor Ostlund, Mater Mahnashi, and Fathi Halaweish. "Estrone Analogs as Potential Inhibitors Targeting EGFR-MAPK Pathway in Non-Small Cell Lung Cancer." Chemical Biology & Drug Design (2023).

[12]. Bray, F., Ferlay, J., Soerjomataram, I., Rebecca, L.S., Torre L.A., Jemal, A. (2018). Global cancer statistics 2018: "GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. CA: A Cancer Journal for The Clinician," 68(6): 394-424.

[13]. Hussain, M., Bird, J.J., Faria, D.R. (2019). "A study on CNN transfer learning for image classification." In UK Workshop on Computational Intelligence, Springer, Cham, pp. 191-202.

[14]. Thakur, Shailesh Kumar, Dhirendra Pratap Singh, and Jaytrilok Choudhary. "Lung cancer identification: a review on detection and classification." Cancer and Metastasis Reviews 39 (2020): 989-998.

[15]. Schenk, Erin L., Tejas Patil, Jose Pacheco, and Paul A. Bunn Jr. "2020 Innovation-Based Optimism for Lung Cancer Outcomes." The Oncologist 26, no. 3 (2021): e454-e472.

[16]. He, Jie, Ni Li, W. Q. Chen, Ning Wu, H. B. Shen, Yu Jiang, Jiang Li, Fei Wang, and J. H. Tian. "China guideline for the screening and early detection of lung cancer (2021, Beijing)." Zhonghua Zhong Liu Za Zhi [chinese Journal of Oncology] 43, no. 3 (2021): 243-268.

[17]. Abdullah, Dakhaz Mustafa, and Nawzat Sadiq Ahmed. "A review of most recent lung cancer detection techniques using machine learning." International Journal of Science and Business 5, no. 3 (2021): 159-173.

[18]. Agarwal, Aman, Kritik Patni, and D. Rajeswari. "Lung cancer detection and classification based on alexnet CNN." In 2021 6th international conference on communication and electronics systems (ICCES), pp. 1390-1397. IEEE, 2021.

[19]. Ramaswamy, Anuradha. "Lung cancer screening: Review and 2021 update." Current Pulmonology Reports 11, no. 1 (2022): 15-28.

[20]. Jacobs, Colin, Arnaud AA Setio, Ernst T. Scholten, Paul K. Gerke, Haimasree Bhattacharya, Firdaus A. M. Hoesein, Monique Brink et al. "Deep learning for lung cancer detection on screening CT scans: results of a large-scale public competition and an observer study with 11 radiologists." Radiology: Artificial Intelligence 3, no. 6 (2021): e210027.