Automated Cancer Detection using Machine Learning and Image Processing

H.M.U.S.S. Samarakoon Department of Information Technology Sri Lanka Institute of Information Technology Malabe, Sri Lanka P.D.S. Fernando Department of Information Technology Sri Lanka Institute of Information Technology Malabe, Sri Lanka B.A.N. Mendis Department of Information Technology Sri Lanka Institute of Information Technology Malabe, Sri Lanka

R.P.P. Kanchana Department of Information Technology Sri Lanka Institute of Information Technology Malabe, Sri Lanka G.W.D.A. Gunarathne Department of Information Technology Sri Lanka Institute of Information Technology Malabe, Sri Lanka L.O. Ruggahakotuwa Department of Information Technology Sri Lanka Institute of Information Technology Malabe, Sri Lanka

ABSTRACT

Cancer is becoming more prevalent around the globe. Even in Sri Lanka, the total cancer rate has doubled in the last 20 years, with a corresponding increase in cancer-related deaths. Cancer is the second leading cause of hospital death. Therefore, a solution to the problem should be an arrangement to reduce time waste, a correct method of directing the patient to detect symptoms, with highly accurate cancer detection, and a better monitoring system. The proposed system is an arrangement that permits and guides a patient to recognize symptoms on their own, directing them to an appropriate healthcare specialist, accurately detecting cancer in its initial stages and monitoring the patient throughout treatment. While cancer detection systems are analyzed, the existing research studies only use one machine learning methodology at a time to diagnose cancer. In the proposed work, Convolutional Neural Network (CNN), Random Forest and XGB Classifier are applied in detecting the presence of breast cancer, brain tumor, skin cancer and lung cancer which outputs results faster with a higher accuracy. This proposed system will be published as a modern cloud-based application which provides better userexperience and ease-of-use.

Keywords

Breast cancer, Brain tumor, Skin Cancer, Lung Cancer, Machine learning, Cancer detection

1. INTRODUCTION

In Sri Lanka, cancer services are predominantly provided by the state sector free of charge to the general public. The implementation of a national cancer policy on cancer prevention and control has resulted in a notable improvement in cancer care supplied throughout the island. Tobacco control and HPV vaccination are two primary preventative methods. The National Cancer Early Detection Center provides screening programs for certain malignancies such as breast, oral, and cervical cancer (CEDC). Every year, however, numerous cases of breast, lung, brain, and skin cancer are detected and treated.

The aim of this projectis to provide an overview of the current cancer statistics in Sri Lanka and propose a system to easily identify cancer in the initial stages with the help of Machine learning and image processing. This will also increase accuracy by using the trained data models to identify symptoms rather than the manual checking done by the doctors. This research will focus on the most common cancer cases, such as breast, lung, skin, and brain tumors. This projectmainly targets medical institutes such as hospitals and medical clinics.

2. BACKGROUND & LITERATURE REVIEW

The burden of cancer across the world is rising. Within the investigation of the past 20 years, the number of incidents of cancer has doubled with the increased of cancer-related mortality. In Sri Lanka, the second common cause of the hospital mortality become as the causes of cancer. This research provides an overview of the status and proposes a system to easily identify cancer in the earlier years with the help of Machine learning and image processing.

2.1. Breast Cancer

For the examination of breast cancer detection, relevant material from multiple sources is referred. The authors have used the most frequently utilized Breast cancer detection methods such as Random Forest, kNN (k-Nearest-Neighbor), Naïve Bayes, Support Vector Machine (SVM) and Bayesian Networks (BN). [1][2]

RFE and SVM are combined in the SVM Classifier method. RFE is a recursive strategy for selecting dataset features based on the least feature value. As a result, in all rounds of SVM-RFE, the incorrect features (lowest weight feature) are removed. SVM performs by picking crucial samples from all classes, known as support vectors, then separating them using these support vectors to generate a linear function that hyperplane to divide the data set into classes. The Random Forest argument is that a single decision tree can provide either a simple or a specific model.[1] Research shows that the ability of RF (Random Forest) to manage data minorities is one of the main reasons for its use in cancer detection. A tumor, for instance, can be identified as benign or malignant, despite the latter class accounting for only 10% of the input data set.

Naive Bayes classifiers are probabilistic classifiers based on the application of Bayes theory. It seems that although it is naive in that it believes that all features are independent of one another, which is rarely the case in real-world circumstances, Naive Bayes has proven to be effective for a wide range of machine learning tasks.[3].A k-nearestneighbor algorithm can be characterized as an algorithm for determining where a data set belongs based on the data sets that surround it. The method is a regression and classification supervised learning strategy[3]. kNN gathers all nearby data points before processing a new data point. Key criteria in determining the distance are attributes with a high degree of variation.

The [3] accuracy of the kNN algorithm is determined to be 95.90 percent, with just one observation misclassified as benign and four observations misclassified as malignant, as shown in the table below, and by the results they have concluded that kNN is better than the Random Forest approach.Research[4]shows an approach of image processing technique, in which used a series of image processing functions to generate a utilized image for machine learning algorithms.

Several studies have attempted on the application of machine learning to the detection and diagnosis of breast cancer, utilizing various methodologies or a combination of algorithms to improve accuracy. Although by looking at several research studies, a conclusion can be made on the limitations of these approaches by referring to previous studies. For example, SVM classifier is ineffective on huge datasets and in high-end computer vision applications. When the training data is not represented, the Naive Bayes Classifier delivers poor results.[5]

2.2. Brain Tumor Cancer

It is impossible to stress the importance of early detection and diagnosis of brain tumors. Brain tumors are frequently diagnosed using computer-aided diagnostic (CAD) techniques in a systematic and specific manner. A brain tumor is a growth of tissue or the central spine that can cause the brain's normal function to be disrupted. In the USA alone, an estimated 83,570 cancer patients have been diagnosed with brain and other Central Nervous System (CNS) tumors in 2021. [6]

According to the International Agency for Research on Cancer (IARC), brain tumors have a 76% mortality rate. [7] It emphasizes the need to identify brain tumors in the earliest stages possible and guide the patients to follow the necessary therapy. With the advancement of technology, it is now possible to use computer-aided systems to detect brain tumors automatically using magnetic resonance imaging (MRI) and computed tomography scans. Machine learning and deep learning approaches such as convolutional neural networks (CNN) have gained popularity among medical research when performing brain tumor detection and categorization. In addition, tumor detection procedures should be performed with exceptional speed and precision.

From complicated medical pictures, suspicious areas are extracted using MR image segmentation. Experts use manual methods to diagnose a brain tumor. Brain tumors are detected by abnormal blobs in the brain during MRI scans (or any other scan). These blobs or areas are brighter than the rest of the brain and have a different lighting than the rest of the brain. The method of segmenting tumors in MRI scans, on the other hand, is extremely challenging. Tumors come in a variety of shapes, sizes, textures, and even locations. If segregating the tumor based on attributes like light, it can run into problems like pixel intensities overlapping with normal tissues. Brain tumor identification and segmentation in MRI images is critical since it reveals the existence of aberrant tissues for therapy or patient follow-up.

2.3. Skin cancer

Cancer is one of the major healthcare burdens across the world. Global statistics suggest almost 10.0 million deaths (9.9 million excluding non-melanoma skin cancer) due to cancer in the year 2020. The most diagnosed cancers include breast cancer in females, lung cancer, and prostate cancers. Lung, liver, and stomach cancers are the major contributors of cancer related deaths. Skin cancer, including both malignant melanoma and non-melanoma skin cancer (NMSC), are common cancers in Caucasians and their incidence is on the rise. According to the US Skin Cancer Foundation, skin cancer affects more people in the United States each year than all other cancers combined [6].

The annual incidence of melanoma cases has increased by 53%. This is due in part to increased ultraviolet (UV) exposure. Even though melanoma is one of the deadliest types of skin cancer, early identification can lead to a high chance of survival [7].

According to the World Cancer Research Fund (WCRF), melanoma of the skin is the 19th most common cancer in both men and women. In 2018 there were 300,000 new cases. Skin cancer, both non-melanoma and melanoma, is the fifth most prevalent kind of cancer in both men and women. In 2018, there were over one million diagnoses globally. According to the American Academy of Dermatology Association (AADA) the most frequent malignancy in the USA is skin cancer. Between 1976-1984 and 2000-2010, the total incidence of BCC (Basal Cell Carcinoma) grew 145%, whereas the overall incidence of SCC (Squamous Cell Carcinoma) climbed 263%. Nonmelanoma BCC and SCC malignancies were more common in women than in males.

According to the American Cancer Association Skin cancers have become increasingly prevalent over the last 30 years, with 87,000 new melanoma diagnoses per year only in the United States.

According to Ferris et al., 2017 with regards to skin melanoma, medical professionals struggle most with the diagnosis rather than the treatment.

2.4. Lung cancer detection

Considering the technological background of this research, can identify that the machine learning models are the main thing of these kind systems and to get better accuracy they should be trained with a larger dataset.

According to some researchers[8]happened in the last few years the referenced research shows how the performance of those research and considering that show some specific technologies which can be used for identifying lung cancers with machine learning. The specific thing that can be identified is, researchers are some of them only have some accuracy likes 90%. Considering that research, the gap can be identified in how they have used those technologies and which kind of datasets they have taken. According to them the research provides an idea about the generated results. Some of them have used only one type of dataset for research.

Considering the nodules of lung cancer and the research along with them[9]which provides an idea about lung cancer identification using CT (Computed Tomography) scanned images. Along with the facts provided using the research, they have used some preprocessing techniques to clear up the datasets for the processing and then they have used those data to train the model. Accuracy of the training process is also matured in the research.

Machine learning or deep learning algorithms are mostly used by those research[10].

Considering the background and the related works, the research shows accuracy around 90% but the research is focusing on archive rather than that.

3. METHODOLOGY

3.1. Breast Cancer Detection

In this study, the research employs a methodology of using two datasets. One data set includes test-result details such as breast tissue diagnosis (benign or malignant), clump thickness, homogeneity of cell size, etc., while the second dataset set includes ultrasound images of breast tissues. In order toperform the Model 1 training,Using the dataset with includes breast tissue features), the major data source for the implementation is a dataset composed of breast cancer report parameters comprising the radius mean, texture mean, perimeter mean, area mean, etc. of the breast tumor. There are various processes involved in extracting clean, diagnosable breast cancer data. The chosen data set should first be investigated and analyzed. The Wisconsin Original dataset was acquired for this study to train the model.

It is compulsory to explore the dataset before feature extraction. The dataset will be pre-processed and investigated initially. After exploring and pre-processing the dataset, selecting the best model to train has been carried out as the third step. Loading the data, eliminating the empty columns from the extracted data, visualizing, and encoding the categorical data are some of the steps that should be completed before splitting the dataset. The data set needs to be split into two portions, 75% train data and 25% test data. Models are trained utilizing Logistic Regression, Decision Tree Classifier, Random Forest Classifier and K-Neighbors Classifier.[8]

When considering the accuracies that each training technique used, the decision tree classifier achieved a training accuracy

of 99.6% out of all the models tested, while the random forest classifier achieved a training accuracy of 99.5%. After running the models on test data and capturing them on a Confusion matrix, the results showed that the Decision Tree Classifier had an accuracy of 95.1%. In comparison, the Random Forest classifier showed 96.5%. Considering the Training and Testing accuracies, it is more suitable for the Random Forest classifier for the model testing.

In order toperform the Model 2 training,the research usesa dataset which includes Ultrasoundimages of breast cancer tissues. This data set includesultrasound images of benign, malignant and normal breast tissues. Utilizing this pertinent data set collection, the model 2 is trained. The data set wasexposed to data pre-processing, feature extraction, and separating the data set intotest and train data during the model's training process.In the process Convolutional Neural Networks (CNNs) algorithm used along with the Keras sequential modelin order toachievehigheraccuracy.

Using the ensemble approach, which joins the accuracies of the two models, it is possible to further boost accuracy after the two models, resulting in a high accuracy.

3.2. Brain Tumor Detection

In this proposed method, the two models will output results faster with higher accuracy to detect brain tumors as early as possible and to avoid further growth using CT and MRI images data sets together with further tests conducted to identify some of the most common symptoms of a brain tumor such as loss of hearing (Acoustic Neuroma) and vision changes (Loss of Vision).

For the first model Convolutional Neural Network (CNN) with LeNet-5 Architecture was applied. The model consisted of 6 layers including Flatten Layer, Conv2D Layer, Dropout Layer, Activation Functions, Dense Layers and Max Pooling Layer. Categorical Cross-entropy with ReLu and Softmax Activation was used as the loss function and activation.



Figure 1:Tumorous and Non-Tumorous Brain MRI (Kaggle)

3.3. Skin Cancer Detection

Since skin cancer has become one of the major healthcare burdens throughout the world, many technical methodologies have been invented to determine and prevent this dangerous cancer. There are Machine Learning based methodologies, but those systems only focus on determining whether the user is diagnosed with skin cancer or not. Other than that, there are no functionalities in those systems. Hence the main goal of implementing this system is to take this cancer detection to

a different level by not only detecting if there is a skin cancer or not but determining exactly what this skin cancer type could be. Because unlike other cancers there are several cancer types in skin cancer, and for each cancer type the symptoms and treatments are different. Hence detecting whether the user has skin cancer or not is not enough. That is why it is so crucial to know the cancer type a patient has diagnosed with when it comes to skin cancer. If this model detects the uploaded image as a skin cancer, then it will predict the three highest probability diagnosed cancer type for that uploaded image of the skin tumor. After that stage there are not any functionalities in all the existing systems. Hence in this system after user gets identified the cancer type there is new functionality for the users for monitoring their cancer status and provide the live status of the cancer to the patients as it grown or shrunk. Monitoring is an important aspect because even if the cancer hasbeen detected the threat is still on if it is not monitored regularly until it is completely cured.

Mainly four main objectives will be covered through this system.

- Suggest users answer some questions to find out how risk they are from diagnosing skin cancer. And provide feedbackand instructions according to the points they get.
- Scan the area of a tumor and determine if it is a skin cancer or not and if it is a cancer what are the three highest probability diagnosed cancer types it could be.
- If the user has a skin cancer, then monitor the diagnosed area of the cancer frequently by comparing the scannedimages and provide the live status of their cancer to the patients as it has grown or shrunk.
- •Separate section to read more about skin cancers, itssymptoms, and methods to avoid it.

To achieve those objectives, a dataset which has some rich data of distinct types of skin cancers was needed. Then train the model from that dataset which brings a higher accuracy. For that the Skin Cancer MNIST: HAM10000 dataset from Kaggle has been used. It contains 10,000 images of sevenseveral types of skin cancers. These seven types of skin cancers are as follows.

Skin Cancer Type	No of Data(Images)
Melanocytic nevi (nv)	[6705 images]
Melanoma (mel)	[1113 images]
Benign keratosis (bkl)	[1099 images]
Basal cell carcinoma (bcc)	[514 images]
Actinic Keratoses (akiec)	[327 images]
Vascular (vasc)	[142 images]
Dermatofibroma (df)	[115 images]

Table 1:Dataset size of each type of skin cancer

The main challenge of using this dataset was the unbalanced data for each cancer type. If the model is built using data as this, it will be biased to one cancer type. And the small amount of data also was a challenge. To overcome those obstacles, Data Augmentation technique has been used. After the data augmentation, the dataset has been expanded to overall nearly 40,000 and nearly 6000 images for each cancer type. A snapshot of the dataset as follows.



Figure 2:Sample of MNIST: HAM10000 dataset from Kaggle

The built model classifies skin lesions into seven classes. The architecture used to build the model was MobileNet CNN. To convert and run this model in the browser the new library called Tensorflow.js has been used. As model Layers ZeroPadding2D Layer, Conv2D Layer, Batch Normalization Layer, ReLU Activation DepthwiseConv2D Layer and Dense Layer have been used. As Loss Function and Optimization Categorical Cross-entropy with Dense and SoftMax Activation and Adam Optimizer have been used.

3.4. Lung Cancer detection

Cancers can be various along with the place where they are situated. The categorized section is differentiated in studying by separating the cancers with those aspects. Lung cancer causes many deaths, which are related to cancers and the main four cancers around the world are lung cancers. According to the detection time of the cancer show, a shocking fact is that the late treatment for cancer can be identified as a major issue of incensement of the death rate. There is another fact that causes lung cancer, which is the high smoking habits and the air contamination around the world.

Investigating the cancer, there are some specific symptoms for lung cancers which signify that with a tumor in their lung.Some of the lung cancer symptoms are as follows,

- Chest pain
- Shortness of breath
- Neck and supraclavicular lymphadenectasis
- Weight loss
- Metastasis pain
- Fatigue
- Fever and dyspnea.

The above symptoms can be considered more common for people to check whether they can be suspected of lung cancer patients. The greater percentage of chest pain, shortness of breath, dyspnea, weight loss, and fatigue can be considered as the symptoms in patients in stage iv, which can be differentiated with other ages along with these symptoms. These symptoms can be various, along with sex and some smoking habits.

This approach for detecting lung cancers using the selected characteristics can be used by building machine learning models and then predicting the probability of having lung cancer or not.

This model would give more accuracy than old models, focusing only on specific data types like numerical or images.

When collecting data for this process, the research data gathering found datasets of both types. The diagram below shows us how the data process is going.



Figure 3 : Methodology of detecting lung cancer

Considering the feature extraction, the identified attributes that can be used to develop the model. After selecting attributes, proceed preprocessing the data to remove the null values and any other error-prone values, reducing the accuracy of the machine learning model. Various classification methods have been used for this numerical dataset. Linear Regression, Random Forest Classifier, K Neighbors Classifier, Decision Tree Classifier, Gradient boosting Classifier XGB Classifier, and Support Vector Classifier are some of them. Considering these classification methods, the best results gained from Random Forest Classifier and XGB classifier. Here XGB will choose because it performs better than the random forest classifier.

Next, thoroughly consider the other model based on the image dataset. Convolutional Neural Networks (CNNs) algorithms used for image classification[9]. These are some other algorithms to be evaluated with this algorithm. According to past research, the VGG-3 version will give more accurate results. Hencethat model will be used for further analysis. Keras sequential model can also be used for this image classification model. It will get good accuracy.

Considering both models and the methodology which have followed, the objective of identifying lung cancers with good accuracy can be achieved using ensemble learning models using these image and numerical models.

When talking about the obstacles whichmust be faced, it shows that here would be some issues with data preprocessing but did solve. The accuracy must be clarified that what is predicted by these models should be evaluated with real data. Predicted data should contrast in a similar perspective to identify the practical ability to use the system in medical environments.

4. RESULTS

4.1. Breast Cancer Detection

Decision Tree Classifier had an accuracy of 95.1%. In contrast, Random Forest Classifier had a 96.5% accuracy rate after the models were run on test data and recorded on a

confusion matrix. It is clear from comparing the Training and Testing accuracies that the Random Forest classifier is a better choice for model testing.

Га	ble	2	:	Testing	Accurac	ies

Algorithm	Accuracy
Logistic Regression	94%
Decision TreeClassifier	95.1%
Random ForestClassifier	96.5%
K-Neighbors Classifier	95%

4.2. Brain Cancer Detection

The following results were obtained using both Binary and Categorical Cross-entropy functions. From the results, the highest accuracy (98%) was obtained by using the Categorical Cross-entropy function with 22 Epochs.

Table 3 : Accuracy obtained by using different loss functions and numbers of epochs

Loss Function	Epochs	Accuracy
Binary Cross-entropy	10	96.3%
Categorical Cross-entropy	20	97.4%
Categorical Cross-entropy	22	98.3%



Figure 4: Accuracy Metrics of Brain Tumor Model



Figure 5: Loss Metrics of Brain Tumor Model

4.3. Skin Cancer Detection

Using Categorical Cross-entropy functions the model was able to obtain a good accuracy rate of 97.8% by the time it reached 30 epochs.

Table 4:Loss function accuracy

Loss Function	Epochs	Accuracy
Categorical Cross-entropy	10	91.7%
Categorical Cross-entropy	20	95.2%
Categorical Cross-entropy	30	97.8%

4.4. Lung Cancer Detection

Below are the results gained by developing machine learning models, and the accuracy of each model is shown in the table. According to this and the dataset research have used as the numerical dataset, research can select the most accurate algorithms are Random Forest Classifier and XGB classifier.

Table 5: Accuracy	obtained	by using	different	algorithms
Table 5. Recuracy	obtaineu	by using	uniterent	angorithmis

Algorithm	Results
K-Nearest Neighbors(KNN)	93.51%
Random Forest Classifier	99.07%
Linear Regression Algorithm	64%
XGB Classifier	99.07%

5. CONCLUSION AND FUTURE RESEARCH

This work is principally centered around advancing predictive models to accomplish great precisionin foreseeing legitimate disease results utilizing supervised machine learning techniques. The analysis of the outcomes connotes that the mix of multidimensional data along with different classification, feature selection, and dimensionality reduction techniques can give favorable tools for inference in this domain. Further research in this field ought to be done to execute the grouping procedures better so it can foresee more factors.

6. REFERENCES

- [1] D. Bazazeh and R. Shubair, "Comparative study of machine learning algorithms for breast cancer detection and diagnosis," 2016.
- [2] A. R. V. B. Soni and S. R. K, "Breast cancer detection by leveraging Machine Learning," vol. 6, no. 4, pp. 320-324, 2020.
- [3] S. Shubham, A. Archit and C. Tanupriya, "Breast Cancer Detection Using Machine Learning Algorithms," in 2018 International Conference on Computational Techniques, Electronics and Mechanical Systems (CTEMS), 2018, 2018.
- [4] A.-H. M. R. A. A. and A. Mohannad, "Breast Cancer Detection Using K-Nearest Neighbor Machine Learning

Algorithm," in 2016 9th International Conference on Developments in eSystems Engineering (DeSE), 2016, 2016.

- [5] S. Kharya and S. Soni, "Weighted Naive Bayes Classifier: A Predictive Model for Breast Cancer Detection," International Journal of Computer Applications, vol. 133, p. 0975 – 8887, 2016.
- [6] K.D.Miller, Q.T.Ostrom, C.Kruchko, N. Patil, T.Tihan, G.Cioffi, H.E.Fuchs, K.A.Waite, A.Jemal, R.L.Siegel, S.Barnholtz, "Brain and other central nervous system tumor statistics," a cancer journal for clinicians, 2021.
- [7] V.S. Lotlikar ,N. Satpute, A. Gupta., "Brain Tumor Detection Using Machine Learning and Deep Learning: A Review. Curr Med Imaging," in Curr Med Imaging, 2022.
- [8] P. Chaturvedi,A. Jhamb, M. Vanani, V. Nemade, "Prediction and Classification of Lung Cancer Using Machine Learning Techniques," in IOP Publishing Ltd, Jaipur, India, 2021.
- [9] S. P. a. H. Z. Rahman, "A New Method for Lung Nodule Detection Using Deep Neural Networks for CT Images," in Int. Conf. on Electrical, Computer and Communication Engineering (ECCE) pp. 1-6, 2019.
- [10] N. M. a. A. L. C. Pehrson, "Automatic pulmonary nodule detection applying deep learning or machine learning algorithms to the LIDC-IDRI database: a systematic review Diagnostics," 2019.
- [11] M. K. Monika, N. A. Vignesh, C. U. Kumari, "Skin cancer detection and classification using machine learning, Materials Today: Proceedings".
- [12] J. Ferlay, M. Colombet, I. Soerjomataram, D.M. Parkin, M. Pineros, A. Znaor, F. Bray, "statistics for the year 2020: An overview.," in Int. J. Cancer, 2021.
- [13] Z. Apalla, D. Nashan, R.B. Weller, X. Castellsaque, "Skin cancer: Epidemiology, disease burden, pathophysiology, diagnosis and therapeutic approaches," in Dermatol. Ther., 2017, 7.
- [14] L.E. Davis, S.C. Shalin, A.J. Tackett, "Current state of melanoma diagnosis and treatment.," in Cancer Biol. Ther, 2019.
- [15] J. Malvehy, G. Pellacani, "Dermoscopy, confocal microscopy and other non-invasive tools for the diagnosis of non-melanoma skin cancers and other skin conditions.," in Acta Derm. Venereol, 2017.
- [16] S. Sengupta, N. Mittal, M. Modi, "Improved skin lesions detection using color space and artificial intelligence techniques.," 2020.
- [17] H.A. Haenssle, C. Fink, R. Schneiderbauer, F. Toberer, T. Buhl, A. Blum, A. Kalloo, A.B.H. Hassen, L. Thomas, A. Enk, "Man against machine: Diagnostic performance of a deep learning convolutional neural network for dermoscopic melanoma recognition in comparison to 58 dermatologists," 2018.