

Skin Lesions Classification and Prediction with Deep CNN

Ms. Ritika Nambiar, Ms. Nisha Sangawar, Ms. Sayali Shinde, Mr. Neeraj Ranade

Department of Computer Engineering

All India Shri Shivaji Memorial Society's Institute of Information Technology, Pune, India

Abstract: *Given the success of Deep Convolutional Neural Network in Computer Vision tasks such as image classification, object detection, etc., DCNN has been applied to many other fields and lays the path for new research domains. Recently, by transfer learning, Esteva et al proposed in "Dermatologist – level classification of Skin Cancer with Deep Neural Networks" that "CNN achieves performance on par with all tested experts, demonstrating an artificial intelligence capable of classifying skin cancer with a level of competence comparable to dermatologists". The results of experiments verify the intuition that features learned by pretrained models and the architectures of the DCNNs help learning features for a completely different domain dataset, here is the skin lesions dermatoscopic images dataset. Given the computational time and the test accuracy of fine-tuning the top layers and fine-tuning the whole model, for this particular dataset, I find that it's better to fine-tune the whole pretrained model with fewer epochs and less computational time and achieve better accuracy.[1].*

Keywords: Convolutional Neural Networks, Transfer Learning, Artificial Intelligence, Deep Learning.

I. INTRODUCTION

Skin cancer is the most widespread cancer diagnosed in the world. It is seen that if it can be diagnosed in its early phases, with choosing the appropriate treatment, survival rates are very good. Hence it is absolutely necessary to get to know at the earliest whether the symptoms of the patient correspond to cancer or not. Traditionally, doctors have been using their naked eye for skin cancer detection. However, on many occasions this leads to not so accurate detection as people make mistakes. Even experts have a tough time saying it, especially when the cancer is at the very early stages. This is where computer vision can help in automating the whole pipeline. There are two types of images available for skin cancer detection. The skin image is captured by the specialized dedicated system in the pathological center with focus on the region of interest with high zoom (E.g., 20x), which needs a skilled dermatologist to conclude the image as positive or negative [5]. This type of image can be fed to a computerized semi-automated system for classification. But in this technology, the victim always needs to walk-into the pathological center and need to take the consultation of a skilled dermatologist [6]. On the other hand, if there is computer software which can automatically detect skin cancer from a digital image captured by any digital image capturing system with little focus on the region of interest, the victim can anytime perform the test even at home. A deep neural network can be trained on thousands of images of both the categories i.e., benign and malignant. By learning the nonlinear interactions, the model can tell whether a new image corresponds to a benign or malignant class. Deep Learning method employs multiple processing layers to learn hierarchical representation of data. It offers a way to harness a large amount of data with few hands feature engineering. Deep Learning method has made impressive advances and evolvment in Computer Vision in recent years, starting from AlexNet in 2012. The image classification task is a very generalized problem: any problems requiring distinguishing between images of different entities can fall within this category. The special characteristics of Deep CNN is that the first layers usually learn very general and "low-level" features of images, while the very last layers of the network learn the semantics and high-level features. Hence, with fine-tuning, Deep Convolutional Neural Networks trained on an image classification task on one dataset can be reused for another image classification task with a different dataset. Consequently, fine-tuning has been used widely in Computer Vision research. Deep neural networks have been found to be very efficient for classification problems among various types of machine learning techniques. Due to the non-linear behaviour of the technique, it can be effectively applied to images as well. Convolutional neural network (CNN) consists of a mass of convolutional modules, and every module usually contains three kinds of layers;

convolutional layer, pooling layer and fully connected layer. CNNs are a kind of neural network which have proven to be very powerful in areas such as image recognition and classification. CNNs can identify faces, pedestrians, traffic signs and other objects better than humans and therefore are used in real time applications like robots and self-driving cars. CNNs are a supervised learning method and are trained using labeled data given with the respective classes. CNNs learn the relationship between the input objects and the class labels and comprise two components: the hidden layers in which the features are extracted and, at the end of the processing, the fully connected layers that are used for the actual classification task. The hidden layers of CNN have a specific architecture consisting of convolutional layers, pooling layers and activation functions for switching the neurons either on or off.

II. LITERATURE SURVEY

Many authors use different machine learning techniques for the classification of skin diseases. Some of the most widely used machine learning techniques like support vector machine artificial neural network decision trees are used by many researchers earlier, but nowadays deep learning-based approach is becoming more popular for the classification of such disease as it gives better classification accuracy and is very suitable for image input. In this extensive literature survey, a traditional method used for classification of skin disease and then deep learning-based approaches used for such classification, furthermore findings of some of the author who uses transfer learning and ensemble learning also discussed in this section. Approach support vector machine and artificial neural network was very common in many after while in deep learning approach convolutional neural network-based model, AlexNet model, ResNet 50/ResNet 101 and dense CNN model are very commonly used by many of the researchers, some researcher also proposed the use of transfer learning where they extracted the feature from deep learning model and apply that feature on some other model. Bashar A. presented a survey on various neural network deep learning techniques with its applications in different areas like speech recognition and image classification. Raj S.J. discussed applications of machine learning methods with its limitation and concluded that the use of deep learning and optimization techniques may increase the accuracy of classification[9]

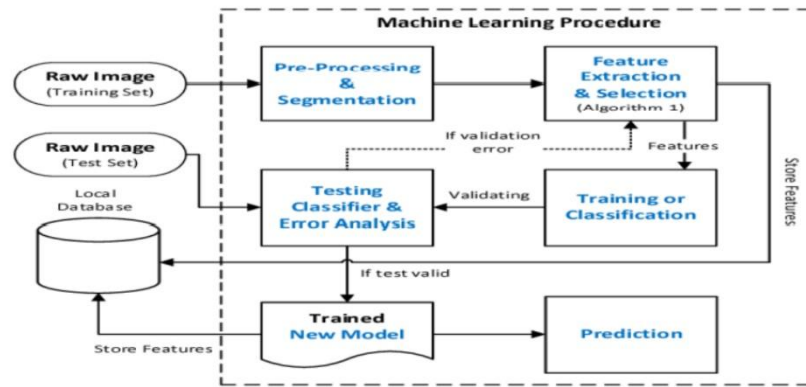
III. RELATED WORK

There has been a lot of work published in the domain of skin cancer classification using deep learning and computer vision techniques. These works use a lot of different approaches including classification only, segmentation and detection, image processing using different types of filters etc. (Esteva et al., 2017) separately used AdaBoost to classify skin lesions. (Xu et al., 2014) used different sets of features including type of lesion, texture, colour etc and neural networks for the making of a diagnosis system. The examples till now only showed algorithms using traditional machine learning techniques, but lately deep learning have proved to be more accurate. The reason is that it automates the feature extraction process completely. It is upto the algorithms to find the better features and train the model accordingly. In (Lopez et al., 2017) made a breakthrough on skin cancer classification by a pre-trained GoogleNet Inception v3 CNN model to classify 129,450 clinical skin cancer images including 3,374 dermatoscopic images. (Rezvantalab et al., 2018) developed an algorithm using Support Vector Machines combined with a deep convolutional neural network approach for the classification of 4 diagnostic categories of clinical skin cancer images. After comparing the performance of number of existing works, it can be stated that CNN architecture is simple to apply for classification of images, but the precision to get the correct images printed in an output is difficult because of the multiple wrong choices of classifiers and pre trained models. Our work builds on the aforementioned approaches. In this paper we have tackled skin cancer classification which is of binary type i.e. there are 2 classes present in the dataset benign and malignant. We have used the publicly available ISIC database for both training and testing the images. We have compared 5 backbone transfer learning architectures - VGG16, ResNet50, InceptionV3, MobileNet and DesneNet169. We have used familiar metrics for evaluating our results including confusion metric. Since there are equal numbers of images in both the classes, hence accuracy alone should be a good enough metric. Nevertheless we present other metric results including precision, recall, F1 score and ROC-AUC. Our work achieves an accuracy of 0.935, precision of 0.94, recall of 0.77, F1 score of 0.85 and ROC- AUC of 0.861 which is better than the previous state of the art approaches.

IV. PROBLEM STATEMENT

Creating a software that can classify skin lesion images and predict the cancer outcomes of the same. The module will take skin images as input and classify them into a suitable category. Then based on the category it will give further predictions regarding cancer detection. We will be training the module using a suitable Dataset. The Dataset will contain images of different skin lesions. Some infected, some non-infected. Then by using machine learning algorithms, the module will get trained. Transfer Learning, Convolutional Neural Networks, Deep Learning, ResNet50 will be the technologies used for the same.

V. PROPOSED SYSTEM ARCHITECTURE



5.1 Image Preprocessing

In the image preprocessing step, we preprocessed our images to enhance the image quality by erasing small pieces of noise for the exact detection of the wanted areas of skin. We get an image as an input and resize the image, then we use median filter to cast off the noise and makes the image noise free. The middle filtering system improves the strength value of the image by exchanging the strength value of the nearby pixel that may have sounds. Then we improved the contrast of the image for finding the better lesion part of the image.

5.2 Feature Extraction

Feature extraction is the preliminary step in image classification. Sometimes the input data size is too large, which is tremendously hard to procedure in its raw form. For solving this, the input data can be transformed into a set of features [5]. After eliminating all the sounds from the image we used Otsu thresholding technique and extract the features of the image. Otsu thresholding is a technique for decrease of a gray level image to a binary image and it includes repeating through all the probable limit esteems and figuring an amount of spread for the pixel levels each side of the edge.

5.3 Classification

To classify data into a given number of classes, we have used machine learning technique. In this paper, the convolutional neural network (CNN) and transfer learning as a feature extractor using 5 classifiers, are used to train the skin lesion images.

VI. CNN ARCHITECTURE

Convolutional Neural Network (CNN) also known as ConvNet is an emerging deep learning algorithm which takes image as input also known as the query image, assigns importance values i.e. biases and learnable weights to different aspects or objects present in the image which makes it possible to differentiate one image from the other. A CNN requires certain pre-processing methods but the best feature of it is that it requires much lower pre-processing with respect to other classification algorithms. While in simple methods, filters are manually applied and are hand-engineered, but with enough training, CNNs acquire the ability to learn from these filters. The main aim of the convolution operation is to withdraw high-level and multi-dimensional features such as edges, colors, shapes, etc. from

the query image. In a particular application of a CNN, there might be one or more number of convolutional layers depending on the complexity of that application. Classically, the first layer is responsible for capturing the lower-level features. As the number-of-layers keep on adding, the architecture to various high-level features giving rise to a dense network. Similar in working to the convolutional layer, the pooling layer is accountable for reducing the spatial size of the convo feature. This is done in order to lower the computational power that is required to process the data obtained through dimensional lowering. The convolutional layer and the pooling layer together, mainly from the layers of a CNN. After going through the above process, the system is now efficient enough to deduce and understand the features of the images.

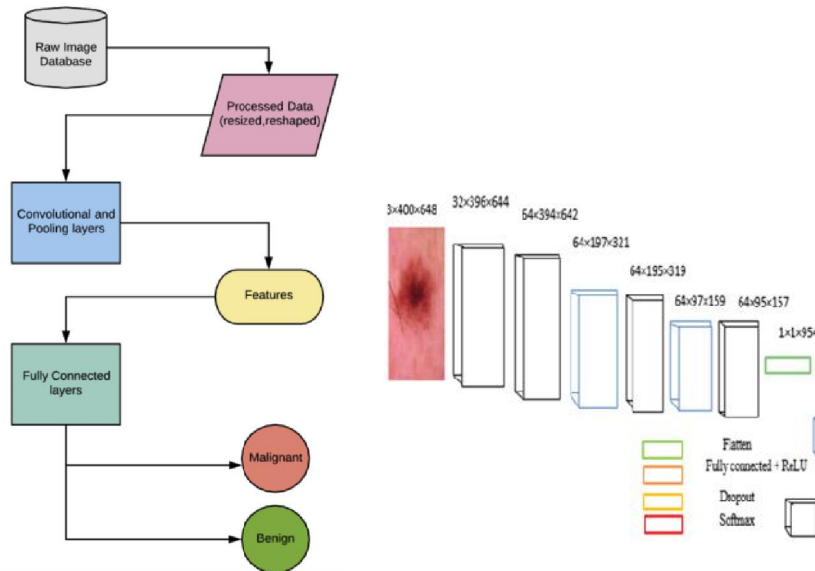


Figure: CNN Architecture

There are three different layers in all CNN architectures:

A) Convolutional Layer

Images are represented as matrices of pixels. In this layer arithmetic operation between image matrix and filter (kernel) matrix is done. Consider image has size $n*n$ and filter has size $m*m$. This filter contains weights for calculation, by this calculation information is extracted from images matrix. Weights in filter combining might be extracting edges, whereas another might focus on colour or they may remove the noise. This filter($m*m$) passes through every pixel in image matrix($n*n$) generating convolutional output matrix. One operation between the image matrix and filter produces one output value which is stored in a convolutional output matrix. There may be many convolutional layers possessing filters for extracting different features. Initial layer extracts basic features, as the network connectivity increases i.e. becomes deeper the next convolutional layer extracts complex and extremely fine features. More clear and proper features from the image and help in better prediction of results. Suppose that the filter matrix moves across the input image matrix by one pixel at a time to cover whole image, then that parameter value one by which filter moves across the image matrix is called stride. While moving filter by stride value over the image matrix the information present at the border of image matrix is not extracted properly as filter moves very less times through it. The output convolution matrix dimensions are less than image matrix dimensions as it depends on filter size. To overcome this both drawbacks extra layers at the border are added to image matrix. This process is called padding. Usually the same values added in padding layers. If zero is taken as values of pixels then it is referred as zero padding. Due to padding feature information at border is extracted properly also the dimension of output convolution matrix is kept same as the input image matrix[6].

B) Pooling Layer

This is an optional layer used for reducing the size of number of parameters when input image size is too large. This layer is added after the convolutional layer and usually added periodically in model. Only motive behind the pooling

layer is to reduce the spatial image. Depth of image remains same after the pooling process. There are different ways for pooling and max pooling is a generally used technique. The output dimensions of images depend on three variable components such as number and size of filter, Size of stride and size of pad added during padding. The output size can be calculated using a simple formula deduced by these variables. The formula is represented

$$\text{as } [(i+f+2p)/s]+1$$

Where, i is input image size, f is size of filter, p is padding size and s is value of stride by which filter is moving.

C) Fully Connected Layer

The previous two layers that are convolutional and pooling layers are only useful for feature extraction and reduction of number of parameters present in image. For generating final output we required a fully connected layer. The training of models based on the features extracted by Convolutional layer is performed by fully connected layers. This layer is same as a normal neural network and possesses loss functions to reduce error in prediction. Similar to normal neural network backpropagation is done for updating weights and biases and for error reduction. After the execution of fully connected layer classification of images based on features extracted is done for application purpose. The process can be explained in detail, flow of the working of the CNN algorithm for feature extraction: The image uploaded will be scanned, and the images having similar features to the image in the database will be extracted and shown as a result. Features of the images database are extracted by using CNN and stored in a file. Like wisely CNN extracts the features of images at run time and these features are then compared one by one with the database images similarity is measured and higher the similarity higher the indexing rank

VII. PROPOSED METHOD

The proposed method has concentrated on recognizing the skin lesion which can help the specialists in making appropriate steps near deal with their patients. In our work, we train our dataset with different CNN classifiers to classify the detected skin lesion area. Depending on the features, the classification is performed using CNN network, transfer learning as a feature extractor using transfer learning classifiers. A number of factors can make skin cancer recognition a challenging task. Some of these factors include flaws in image quality like uneven brightness, obstruction, and also the fact that there are many images which have similar shape, color and texture. All these factors have been considered by us through the dataset which have chosen and exhibits the correct precision of images.

A) Dataset

We obtained a public dataset from ISIC website for skin cancer classification. We used the ISIC 2018 Challenge Archive Downloader to download the images. We used 3000 images for training and 600 images for validation of size 224×224 . The images are distributed equally between training and validation sets which are shown below in figure.

B) Methodology

We have used the concept of transfer learning for the classification. With transfer learning, instead of starting the learning process from scratch, the model starts from patterns that have been learned when solving a different problem. This way the model leverages previous learnings and avoids starting from scratch. In image classification, transfer learning is usually expressed through the use of pre-trained models. A pre-trained model is a model that was trained on a large benchmark dataset to solve a problem similar to the one that we want to solve. We used three pre-trained models Inception v3, InceptionResNet v2 and ResNet152 as the pre trained weights for our work.

1) Inception V3: Google's Inception v3 architecture was re-trained on our dataset by fine-tuning across all layers and replacing top layers with one average pooling, two fully connected and finally the softmax layer allowing to classify 2 diagnostic categories. The size of input images was all resized to $(224, 224)$ to be compatible with this model. Learning rate was set to 0.0001 and Adam was used for the optimizer.

2) InceptionResNet v2: InceptionResNet v2 architecture was re-trained on our dataset by fine-tuning across all layers and replacing top layers with one global average pooling, one fully connected and finally the softmax layer allowing to classify 2 diagnostic categories. The size of input images was all resized to $(224, 224)$ to be compatible with this model. Learning rate was set to 0.0001 and Adam was used for the optimizer.

3) MobileNet: The fundamental part of MobileNet is depthwise separable filters, named as Depthwise Separable Convolution. These convolutions layers which is a form of factorized convolutional factorize a standard convolution into a depthwise convolution called a pointwise convolution. In MobileNet, the depthwise convolution applies a single filter to each input channel. The pointwise convolution then applies the convolution operation to combine the outputs of the depthwise convolution.

4) DenseNet169: To solve the vanishing gradient problem, this architecture uses a simple connectivity pattern to ensure the maximum flow of information between layers both in forward and backward computation. The layers are connected in a way such that inputs from all preceding layers passes through its own feature-maps to all subsequent layers. To facilitate the down-sampling in the architecture, the entire architecture is divided into multiple densely connected blocks. The layers between these dense blocks are transition layers which perform convolution and pooling operations.

5) ResNet50: ResNet50 architecture was re-trained on our dataset by fine-tuning across all layers and replacing top layers with one average pooling, one fully connected and finally the softmax layer allowing to classify 2 diagnostic categories. The size of input images was all resized to (224, 224) to be compatible with this model. Learning rate was set to 0.0001 and Adam was used for the optimizer. It uses identity mapping to map the inputs. This identity mapping does not have any parameters and is just there to add the output from the previous layer to the layer ahead. The identity mapping is multiplied by a linear projection to expand the channels of shortcut to match the residual. The Skip Connections between layers add the outputs from previous layers to the outputs of stacked layers. This results in the ability to train much deeper networks than what was previously possible. But trying all these pre trained models we find that REsNet50 gives us the more accurate answers than the latter pre trained models.

6) Transfer Learning: Transfer learning is a popular method in computer vision that allows us to build accurate models faster. With transfer learning, instead of starting the learning process from scratch, we start from patterns that have been learned when solving a different problem. The advantages of using transfer learning are: Super simple to incorporate. Achieve the same or even better (depending on the dataset model performance quickly. There's not as much labeled data required. Versatile uses cases from transfer learning, prediction, and feature extraction. The proposed method can be summarized in the following seven points:

1. We split the dataset into two parts-training set and test set with 80 percent and 20 percent images respectively.
2. We used data augmentation like shearing, zooming, flipping and brightness change to increase the dataset size to almost double the original dataset size.
3. We tried with pre trained models like Inception v3, InceptionResNet v2, ResNet 50, MobileNet and DenseNet169 by fine tuning the last few layers of the network.
4. We used 50 percent dropout and batch normalization layers in between to reduce overfitting.
5. We used two dense layers with 64 neurons and 2 neurons respectively. The last layer is used for the classification with softmax as the activation function.
6. We used binary cross entropy as the loss function. 4
7. We trained the model for 20 epochs with a batch size of 32 by changing the hyper-parameters like learning rate, batch size, optimizer and pre-trained weights.

We see that all these transfer learning methods come under one category where they are trained, having some instances same, in terms of the architecture. Input image is given to the 1st layer of convolutional layer of VGG net. Images are represented as matrices of pixels of height weight as 224 x 224 and accept 3 channels that are RGB – Red, Green and blue colored images as input. These convolutional layers are used as filters whose parameters are to be acquired though learning. Convolutional layers along with Rectified Linear Unit (ReLU) activation function are used to learn new features with less similarity and error rate is also minimized. ReLU function is represented as: Here, x is the input. The convolutional layers extracts fine, complex and clear features which help in better prediction. The output convolutional matrix has size less than then image matrix which depends on the filter the paddings(extra layers added to the image matrix) are added to extract features from the borders properly.



The output of convolutional layers is then passed to max pooling layers. Max pooling reduces the number of parameters and takes only prominent features and is used to scale in variant features and helps reducing over-fitting. This reduces the size of matrix by dividing the matrix in parts and taking the maximum value.

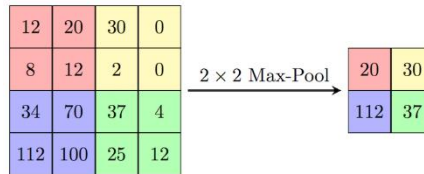


Figure: Max Pooling on a matrix

The output from the convolutional and pooling layers is in the form of matrix which is then converted to vector by flatten layer and is given to fully connected layer. Fully connected layer is a feed forward neural network in which information moves only in forward direction and is used for classification. These layers are similar to artificial neural network and they perform similar computational operations. The model processes these images and gives output as a vector. This vector consists of classification probabilities. The softmax function is used at the end so that all these probabilities add up-to 1.

Output vector can be represented by:

$$\text{Softmax}(x_i) = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$

Here N is the number of classes, x is the input vector. The features of the dataset are calculated and stored in a file then at run time the features of query image are calculated. Dot product or inner product is used for measuring the similarity between the image vector uploaded and image vectors from the database. This helps to find the set of images similar to the uploaded image. Efficiency is calculated using mean average precision. It is given by

$$\text{Precision} = \frac{TP}{\text{total positive results}}$$

$$\text{Recall} = \frac{TP}{\text{total cancer cases}}$$

VIII. RESULT ANALYSIS AND DISCUSSION

In this section we present our findings. We plotted the loss vs epochs, accuracy vs epochs, confusion matrix for the classifier and ROC-AUC curve for the classifier.

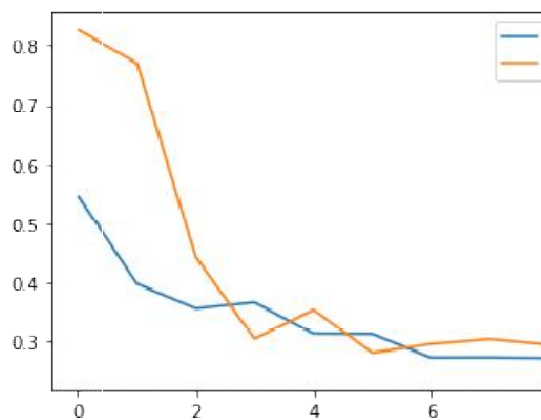


Figure: Loss v/s epoch

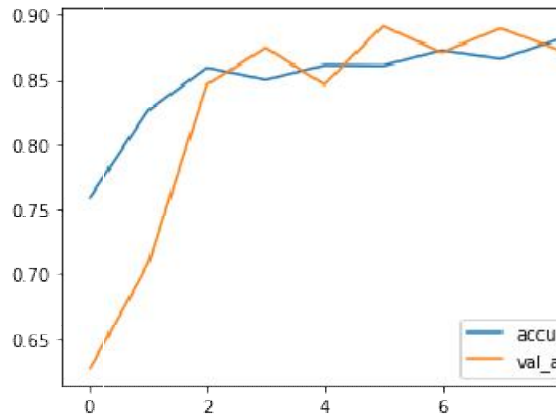


Figure: Accuracy v/s epoch

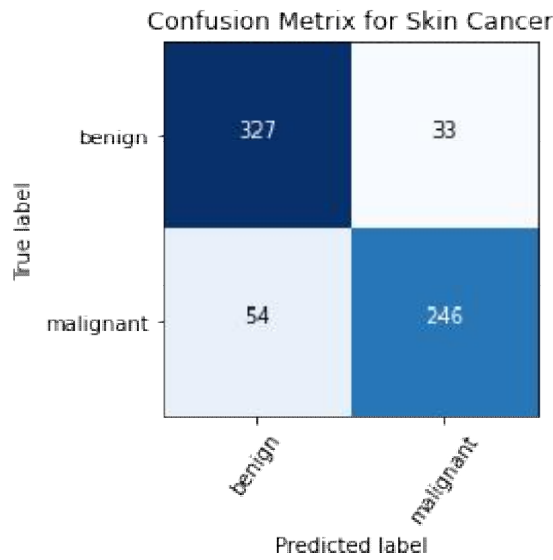


Figure: Confusion Matrix

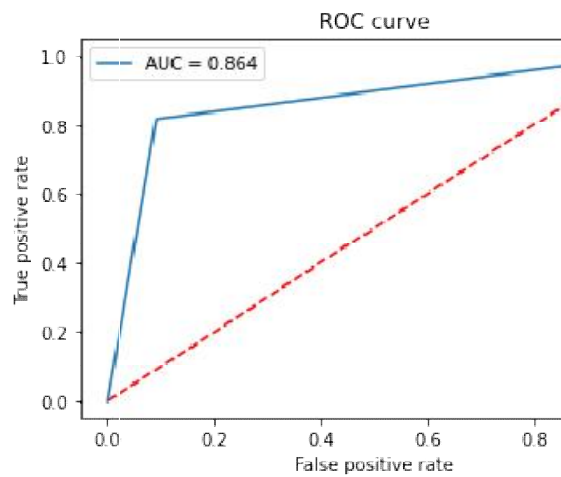
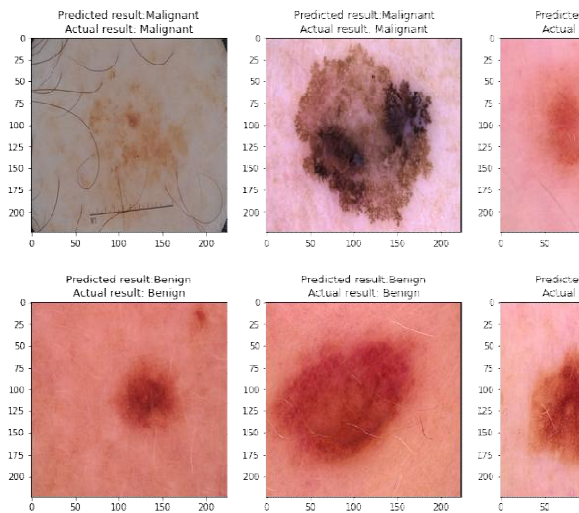


Figure: ROC-AUC

And also right predictions have been achieved through the enormous processes which we have done. It gets segmented by predicted results and actual result .The following figure is below:



The details of the experiment demonstrated the effect of training dataset size is shown in Table 1.

Train Size	Precision	Recall	F1Score
300	0.87	0.74	0.85
600	0.80	0.86	0.86

Table 1: Precision, recall and F1 Score values

Next, we present our findings and show the validation of our trained models. In this paper, two types of major skin cancer categories are used. The evaluation and results of trained models is calculated by common classification metrics. The ROC curve is calculated by plotting the sensitivity against 1-specificity and can be used to evaluate the classifier. The further the ROC curve deviates from the diagonal, the better the classifier. We found that a batch size value of 64 gives better results when compared to that of 32 as shown in Table 2.

Batch Size	Accuracy	Precision	Recall F1	Score	ROC-AUC
32	0.885	0.87	0.74	0.82	0.844
64	0.892	0.88	0.82	0.85	0.864

Table 2: Effect of batch size on results

Pre Trained	Accuracy	Precision	Recall	F1 Score	ROC-AUC
ResNet50	0.892	0.88	0.82	0.85	0.864
InceptionV3	0.88	0.87	0.81	0.84	0.839
Inception ResNet v2	0.873	0.86	0.80	0.82	0.829
MobileNet	0.854	0.842	0.78	0.81	0.824
Densenet169	0.836	0.82	0.76	0.80	0.817

Table 3: Comparison of pre trained weights on results

IX. CONCLUSION AND FUTURE WORK

In conclusion, this study investigated the ability of deep convolutional neural networks in the classification of benign vs malignant skin cancer. Our results show that state-of-the-art deep learning architectures trained on dermoscopy images (3600 in total composed of 3000 training and 600 validation) outperforms dermatologists. We showed that with use of very deep convolutional neural networks using transfer learning and fine-tuning them on dermoscopy images, better diagnostic accuracy can be achieved compared to expert physicians and clinicians. Although no preprocessing step is applied in this paper, the experimental results are very promising. These models can be easily implemented in dermoscopy systems or even on smartphones in order to assist dermatologists. More diverse datasets (varied categories, different ages) with much more dermoscopy images and balanced samples per class is needed for further improvement. Also using the metadata of each image can be useful to increase the accuracy of the model. The overall accuracy and loss of the network indicate satisfactory outcome that can be improved further. The next steps in this research include covering a wider variety of skin cancer types.

REFERENCES

- [1]. Xu, J. S. Ren, C. Liu, and J. Jia. Deep convolutional neural network for image deconvolution. In Advances in neural information processing systems, pages 1790–1798, 2014
- [2]. A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, and S. Thrun. Dermatologistlevel classification of skin cancer with deep neural networks. *nature*, 542(7639):115–118, 2017.
- [3]. A. R. Lopez, X. Giro-i Nieto, J. Burdick, and O. Marques. Skin lesion classification from dermoscopic images using deep learning techniques. In 2017 13th IASTED international conference on biomedical engineering (BioMed), pages 49–54. IEEE, 2017.
- [4]. A. Rezvantalab, H. Safigholi, and S. Karimijeshni. Dermatologist level dermoscopy skin cancer classification using different deep learning convolutional neural networks algorithms. *arXiv preprint arXiv:1810.10348*, 2018.
- [5]. X. Yuan, Z. Yang, G. Zouridakis, and N. Mullani, “SVM-based texture classification and application to early melanoma detection,” in Engineering in Medicine and Biology Society (EMBS), 28th Annual International Conference on. IEEE, pp. 4775-4778, August 2006
- [6]. MutasemAlsmadi, “An efficient similarity measure for Content Based Image Retrieval using memetic algorithm”, Taylor and Francis, 2019
- [7]. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 4700–4708, 2017.
- [8]. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In Advances in neural information processing systems, pages 1097–1105, 2012.