# Skin Disease Detection

## Final Report

### Supervised by : Dr. Mohamed Elshenawy

---

**Team**

1- Ibrahim Abdalaal 201800224

2- Amr Ragab Abdelaziz 201801041

3- Mohamed Abdelwahed 201801978

4- Abdullah Kamal Muhammad 201801271

## Problem Definition and Motivation

Deep learning (DL),machine learning (ML) and Artificial intelligence (AI) have added many improvements to different medical fields. Due to these improvements, it is now possible to detect diseases earlier and get better results by integrating the efforts of human doctors to the added value of the DL models. There is a medical field called Dermatology which studies, detects and treats skin diseases of skin, hair and nails. Unfortunately,  Dermatology did not achieve great progress like other medical fields [1]. There are many reasons that prevent that medical field from being improved, such as weather patterns and personal hygiene-related medical treatment. The  variety and challenges that face medical professionals to recognise the features of skin disorders, the majority of which are skin diseases.

The use of AI systems will represent a strong addition to the early detection of diseases, which increases the possibility of recovery before it reaches its final stages that traditional medicine cannot treat, as some diseases have developed into cancerous cells [2]. Late detection of diseases increases the risk because of its high cost, and therefore the solution lies in the use of artificial intelligence applications such as computer vision, which will be available to anyone regardless of the financial level of the patient, and all that is needed is a picture of the area believed to have a disease instead of laboratory examinations. DL models will provide us with the needed flexibility and high accuracy of detection we aim to. DL models improve their accuracies by using datasets with large size so it will be more helpful to benefit from the concept of Transfer Learning and use other or pre-trained models that learned from other data that is not available to me which will result in better accuracies.

## Dataset

We used two datasets to improve the training of our model. The first dataset has images for 23 different skin diseases. It has 19500 RGB images. The dataset is split into 15,500 images in the training and 4000 in the test. The images' resolution is not very high. Additionally, the 23 classes are not equally distributed. The dataset is retrieved from the public portal Dermnet (http://www.dermnet.com), which is the largest dermatology source online built for the purpose of providing online medical education.

The second data set is retrieved from the kaggle website and it is called 'Skin diseases image dataset'. This dataset has 10 classes for different skin diseases. This dataset has 27053 RGB images. The dataset is not splitted. The images in it have a better quality than the Dermnet dataset.

We used the second dataset to increase the number of images in the first dataset, so we merged the second dataset with the training in the first dataset to have more quality images and better distribution.

## Approach

We used transfer learning as our main approach.by training three different models on our dataset-VGG16 , ResNet50 and Xception- and comparing the performance of the three models to settle on the best model.

Since these pretrained models are already tested on image classification tasks and taking into consideration the difference between our dataset and the original datasets the models are trained on ,We believe that good accuracy could be achieved by fine training and tuning.

## Methodology

To achieve good accuracy ,these models should have enough data since they are complex ,hence ,data-hungry.With the small size of our dataset along with the hard classification task-23 classes- two solution were proposed as follows :

- Using data augmentation
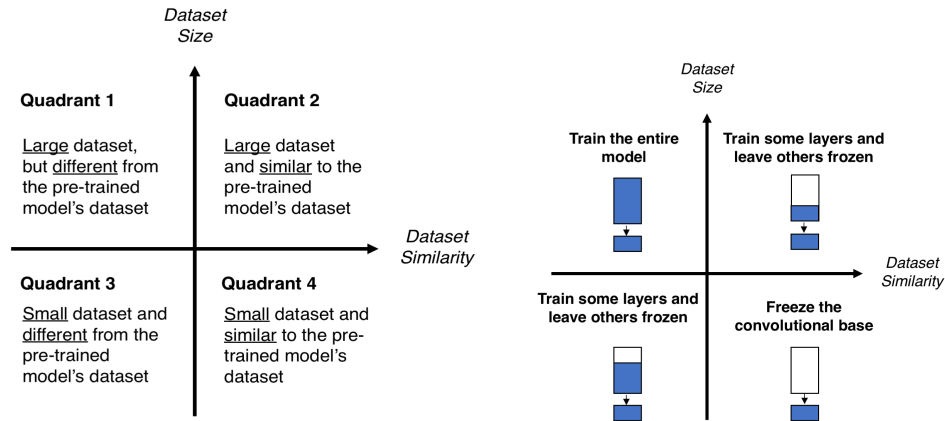- Increasing the data by merging similar dataset with the same  classes if possible

Dataset augmentation gave us two advantages: simulating having a larger dataset and avoiding the overfitting problem that we expected to face.

For dataset merging ,we were able to find another dataset with the same data as our original dataset.By merging the two together (this merging done on the training data only leaving the original test set as it is ) ,we increased the dataset size which increased our models performance.

Next step is to train our models the processes can be summarized as as follows :

1.Selecting the pretrained models to use which are ResNet ,VGG16 and Xception

2.Classifing the problem based on the size similarity matrix.for our dataset which is small in size and different from the pretrained model data set ,Our problem falls within quadrant 3 in the figure below.



3.based on the above step ,we will train some layers leaving the rest of layers frozen

4.Building our classifier which composed of 3 fully connected layer followed by softmax activation layer

5.Dataset augmentation has been done using Keras preprocessing Image data generator for each model.We applied the following data augmentation techniques

- width_shift_range=0.4
- height_shift_range=0.3
- rotation_range=90

We tried different values for each of them to get the best accuracy

6.Apply Zooming range=.2  to the image as we noticed that most of the information about the disease is in the middle of the images with a lot of useless information in the rest of the image.

7-As expected and due to the model complexity and the size of our dataset ,overfitting happens so we used two combined techniques from L2 Regularization and Dropout to reduce its effect.

8.fine Tuning our model with the following hyperparameters to tune:

- Number of frozen to active layers
- Learning rate

- Number of layers and neurons of our classifier
- Dropout rate
- L2 regularization penalty

for Number of frozen layers ,we have tried to unfrozen different numbers of layers starting with 10 layers to 40 layer in the very deep model xception

- learning rate trying .01 ,.001 .0001
- Number of layers our classifier ,trying 5 ,4,3 connected layers
- Number of neurons of our classifier  tying decreasing number with each layer such as 1024 ,512,256
- Dropout rate  : .2 , .5 , .7
- L2 regularization penalty .01 ,.1 , .5

9.Evaluting our model which done by evaluating the accuracy on the test set and by classification report and confusion matrix

## Results

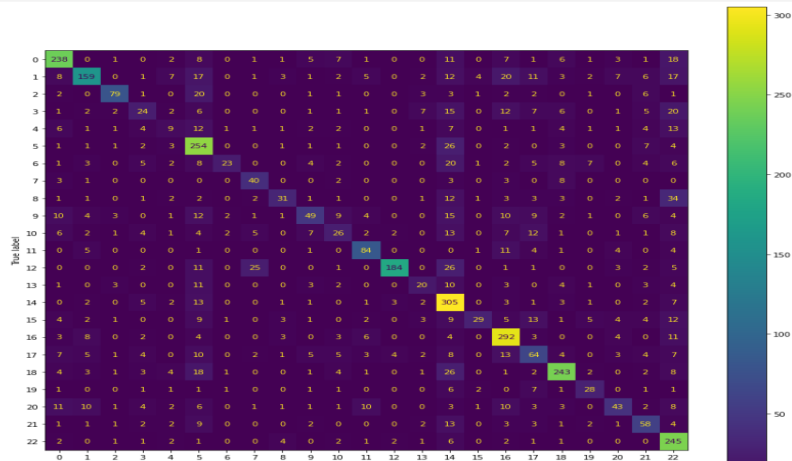We tried multiple tuning as mentioned above in the methodology.

1. **Xception model**
- First, we have model_Xception with model base layers from 30 to 60 are trainable with dropout rate (0.7) and learning rate .0001
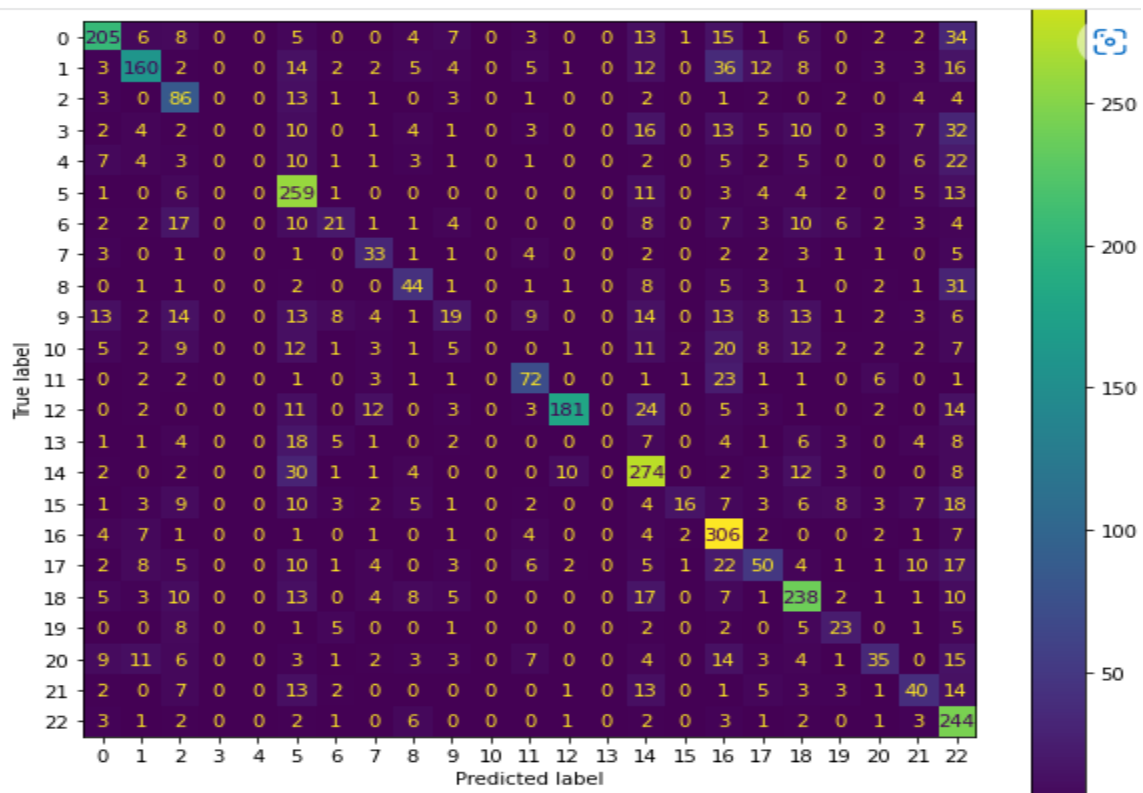
Test accuracy =.6314

```
model_xception.evaluate(test_generator)

63/63 [==============================] - 64s 1s/step - loss: 1.5449 - accuracy: 0.6314
]: [1.544927716255188, 0.6314342617988586]
```

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.77 | 0.76 | 0.76 | 312 |
| 1 | 0.76 | 0.55 | 0.64 | 288 |
| 2 | 0.81 | 0.64 | 0.72 | 123 |
| 3 | 0.36 | 0.21 | 0.27 | 113 |
| 4 | 0.21 | 0.12 | 0.16 | 73 |
| 5 | 0.58 | 0.82 | 0.68 | 309 |
| 6 | 0.74 | 0.23 | 0.35 | 101 |
| 7 | 0.51 | 0.67 | 0.58 | 60 |
| 8 | 0.61 | 0.30 | 0.41 | 102 |
| 9 | 0.56 | 0.34 | 0.42 | 143 |
| 10 | 0.36 | 0.25 | 0.29 | 105 |
| 11 | 0.69 | 0.72 | 0.71 | 116 |
| 12 | 0.94 | 0.70 | 0.81 | 261 |
| 13 | 0.43 | 0.31 | 0.36 | 65 |
| 14 | 0.55 | 0.87 | 0.67 | 352 |
| 15 | 0.72 | 0.27 | 0.39 | 108 |
| 16 | 0.71 | 0.85 | 0.77 | 343 |
| 17 | 0.42 | 0.42 | 0.42 | 152 |
| 18 | 0.79 | 0.75 | 0.77 | 325 |
| 19 | 0.54 | 0.53 | 0.53 | 53 |
| 20 | 0.56 | 0.36 | 0.43 | 121 |
| 21 | 0.49 | 0.55 | 0.52 | 105 |
| 22 | 0.56 | 0.90 | 0.69 | 272 |
| | | | | |
| accuracy | | | 0.63 | 4002 |
| macro avg | 0.59 | 0.53 | 0.54 | 4002 |
| weighted avg | 0.64 | 0.63 | 0.62 | 4002 |

- Xception model with L2 penalty =.001 and dropout rate 0.5 (rest of parameters are the same as above),model base layers from 30 to 60 are trainable

  ->Test accuracy =.5762

  ->Confusion Matrix :

Classification report :

```
              precision    recall  f1-score   support

           0       0.75      0.66      0.70       312
           1       0.73      0.56      0.63       288
           2       0.42      0.70      0.52       123
           3       0.00      0.00      0.00       113
           4       0.00      0.00      0.00        73
           5       0.56      0.84      0.67       309
           6       0.39      0.21      0.27       101
           7       0.43      0.55      0.49        60
           8       0.48      0.43      0.46       102
           9       0.29      0.13      0.18       143
          10       0.00      0.00      0.00       105
          11       0.60      0.62      0.61       116
          12       0.91      0.69      0.79       261
          13       0.00      0.00      0.00        65
          14       0.60      0.78      0.68       352
          15       0.70      0.15      0.24       108
          16       0.59      0.89      0.71       343
          17       0.41      0.33      0.36       152
          18       0.67      0.73      0.70       325
          19       0.40      0.43      0.41        53
          20       0.51      0.29      0.37       121
          21       0.39      0.38      0.38       105
          22       0.46      0.90      0.60       272

    accuracy                           0.58      4002
   macro avg       0.45      0.45      0.43      4002
weighted avg       0.54      0.58      0.54      4002
```

- Second, L2 Regularization (penalty= .01) added to xception model with the same parameters as above.
    - Test accuracy ~0.5480

**Comments on Xception model Performance :**

- Evaluation the model with F1 score ,precision and recall is necessary to get the true performance of our model.
- With L2 regularization combined with dropout ,the overfitting problem is less obvious but still the test accuracy does not increase
- With xception being very deep complex model combined with our small data set ,L2 doesn't seem to do much to the test accuracy as the overfitting problem still exist.
- In terms of these test accuracy , F1 ,precision and recall scores.we could say xception model provide good performance compared to state of the art
- We tried L2 penalties .01 ,.1,.5 but none of them give good accuracy so we settled on using only dropout
- So L2 with dropout seems to over penalize the model ,final decision to use only Dropout with 3 layers in classifier

## 2. VGG16 model

- we have VGG16 with model base layers from 8 to 14 are trainable with dropout and learning rate .0001
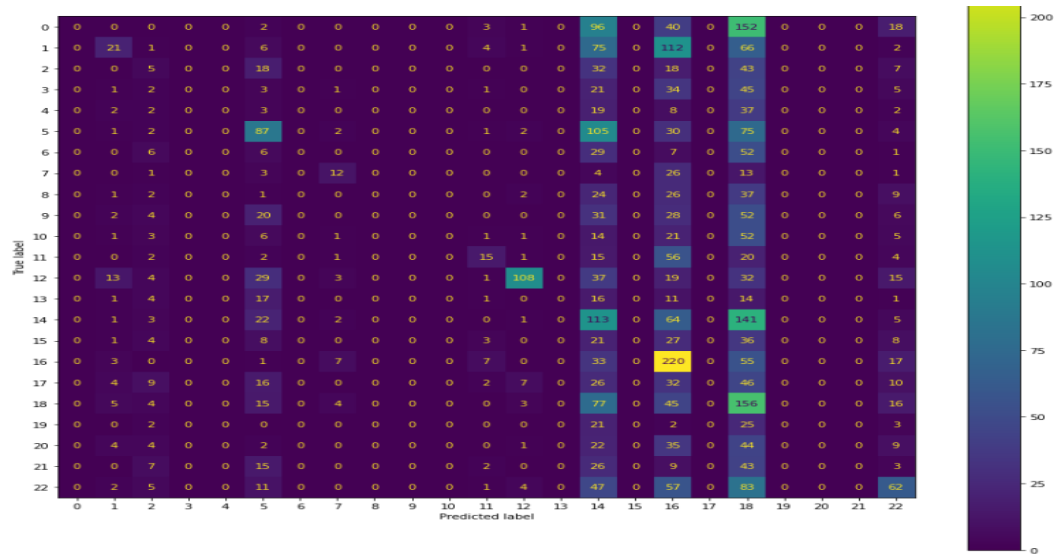
```
Epoch 1/10
409/409 [==============================] - 870s 2s/step - loss: 0.1733 - accuracy: 0.2385 - val_loss: 0.1494 - val_accuracy: 0.3237
Epoch 2/10
409/409 [==============================] - 745s 2s/step - loss: 0.1344 - accuracy: 0.3394 - val_loss: 0.1310 - val_accuracy: 0.3565
Epoch 3/10
409/409 [==============================] - 739s 2s/step - loss: 0.1295 - accuracy: 0.3645 - val_loss: 0.1378 - val_accuracy: 0.3360
Epoch 4/10
409/409 [==============================] - 749s 2s/step - loss: 0.1290 - accuracy: 0.3705 - val_loss: 0.1349 - val_accuracy: 0.3561
Epoch 5/10
409/409 [==============================] - 757s 2s/step - loss: 0.1285 - accuracy: 0.3761 - val_loss: 0.1469 - val_accuracy: 0.3279
Epoch 6/10
409/409 [==============================] - 755s 2s/step - loss: 0.1294 - accuracy: 0.3763 - val_loss: 0.1335 - val_accuracy: 0.3617
Epoch 7/10
409/409 [==============================] - 758s 2s/step - loss: 0.1283 - accuracy: 0.3759 - val_loss: 0.1544 - val_accuracy: 0.3458
Epoch 8/10
409/409 [==============================] - 751s 2s/step - loss: 0.1330 - accuracy: 0.3730 - val_loss: 0.1415 - val_accuracy: 0.3654
Epoch 9/10
409/409 [==============================] - 756s 2s/step - loss: 0.1290 - accuracy: 0.3766 - val_loss: 0.1570 - val_accuracy: 0.3620
Epoch 10/10
409/409 [==============================] - 744s 2s/step - loss: 0.1312 - accuracy: 0.3781 - val_loss: 0.1562 - val_accuracy: 0.3286
```

```
model.evaluate(test_generator)
```

```
63/63 [==============================] - 58s 920ms/step - loss: 0.1640 - accuracy: 0.1997
20]: [0.164026141166687, 0.19965016841888428]
```

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.00 | 0.00 | 0.00 | 312 |
| 1 | 0.33 | 0.07 | 0.12 | 288 |
| 2 | 0.07 | 0.04 | 0.05 | 123 |
| 3 | 0.00 | 0.00 | 0.00 | 113 |
| 4 | 0.00 | 0.00 | 0.00 | 73 |
| 5 | 0.30 | 0.28 | 0.29 | 309 |
| 6 | 0.00 | 0.00 | 0.00 | 101 |
| 7 | 0.36 | 0.20 | 0.26 | 60 |
| 8 | 0.00 | 0.00 | 0.00 | 102 |
| 9 | 0.00 | 0.00 | 0.00 | 143 |
| 10 | 0.00 | 0.00 | 0.00 | 105 |
| 11 | 0.36 | 0.13 | 0.19 | 116 |
| 12 | 0.82 | 0.41 | 0.55 | 261 |
| 13 | 0.00 | 0.00 | 0.00 | 65 |
| 14 | 0.12 | 0.32 | 0.18 | 352 |
| 15 | 0.00 | 0.00 | 0.00 | 108 |
| 16 | 0.24 | 0.64 | 0.35 | 343 |
| 17 | 0.00 | 0.00 | 0.00 | 152 |
| 18 | 0.12 | 0.48 | 0.19 | 325 |
| 19 | 0.00 | 0.00 | 0.00 | 53 |
| 20 | 0.00 | 0.00 | 0.00 | 121 |
| 21 | 0.00 | 0.00 | 0.00 | 105 |
| 22 | 0.29 | 0.23 | 0.26 | 272 |
| | | | | |
| accuracy | | | 0.20 | 4002 |
| macro avg | 0.13 | 0.12 | 0.11 | 4002 |
| weighted avg | 0.18 | 0.20 | 0.16 | 4002 |

**Comments on VGG16**

- The performance of the Vgg16 model is not good with very low test accuracy and f1 score.
- Since VGG16 considered shallow if compared to Xception model ,the performance is expected compared to Xception and the following ResNet Model
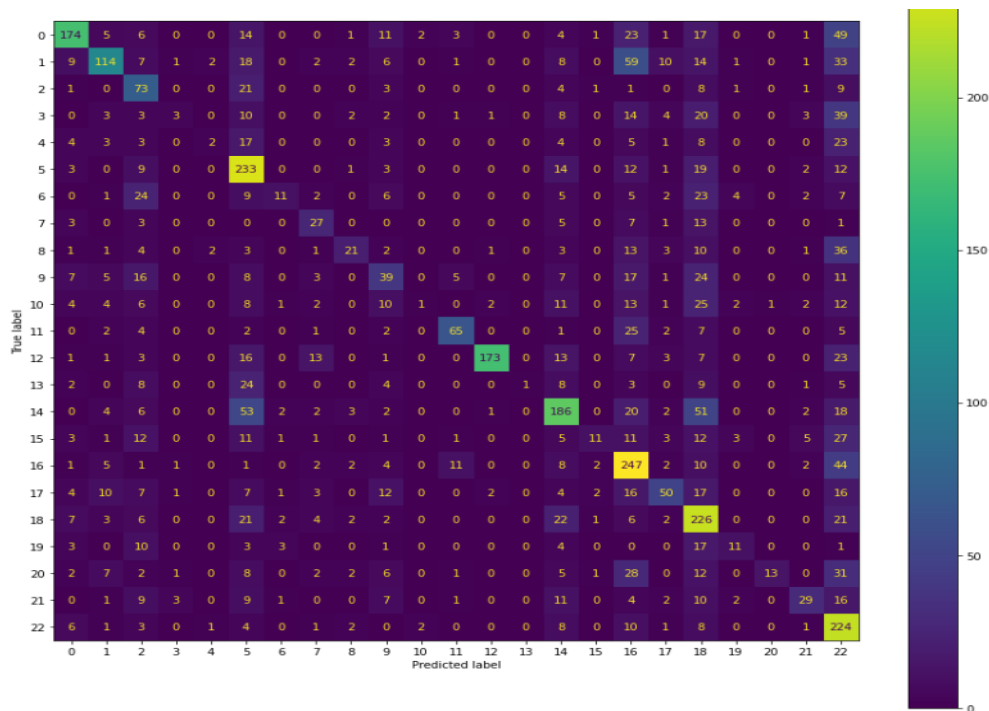
- The main advantage of VGG16 model is less overfitting than Xception and Resnet.This could be justified also due to the vert less number of layers in VGG6
- Not expecting much from VGG for the above reasons

3. **ResNet50**

- we have Resnet50 with model base layers from 1 to 25are trainable with dropout rate by 0.7 and learning rate .0001

```
model_Res50.evaluate(test_generator)
```

```
63/63 [==============================] - 57s 901ms/step - loss: 2.0785 - accuracy: 0.4833
[2.0785465240478516, 0.48325836658477783]
```

```
              precision    recall  f1-score   support

          0       0.74      0.56      0.64       312
          1       0.67      0.40      0.50       288
          2       0.32      0.59      0.42       123
          3       0.30      0.03      0.05       113
          4       0.29      0.03      0.05        73
          5       0.47      0.75      0.58       309
          6       0.50      0.11      0.18       101
          7       0.41      0.45      0.43        60
          8       0.55      0.21      0.30       102
          9       0.31      0.27      0.29       143
         10       0.20      0.01      0.02       105
         11       0.73      0.56      0.63       116
         12       0.96      0.66      0.78       261
         13       1.00      0.02      0.03        65
         14       0.53      0.53      0.53       352
         15       0.58      0.10      0.17       108
         16       0.45      0.72      0.56       343
         17       0.54      0.33      0.41       152
         18       0.40      0.70      0.51       325
         19       0.46      0.21      0.29        53
         20       0.93      0.11      0.19       121
         21       0.55      0.28      0.37       105
         22       0.34      0.82      0.48       272

   accuracy                           0.48      4002
  macro avg       0.53      0.37      0.36      4002
weighted avg       0.54      0.48      0.45      4002
```
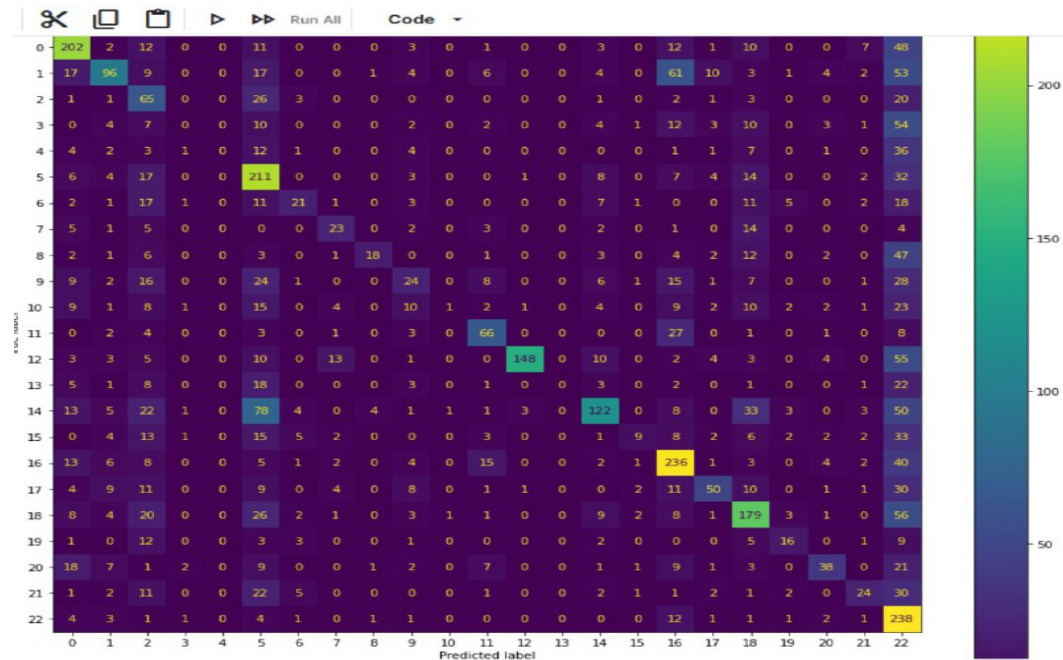
**ResNet50 (another trial)**

- we have Resnet50 with model base layers from 25 to 45 are trainable with dropout rate 0.5 and learning rate .0001

  Test accuracy=.4465

```
model_Res50.evaluate(test_generator)

63/63 [==============================] - 57s 911ms/step - loss: 2.0930 - accuracy: 0.4465
[11]: [2.0929923057556152, 0.44652673602104187]
```

Confusion matrix

Classification Report

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.62 | 0.65 | 0.63 | 312 |
| 1 | 0.60 | 0.33 | 0.43 | 288 |
| 2 | 0.23 | 0.53 | 0.32 | 123 |
| 3 | 0.00 | 0.00 | 0.00 | 113 |
| 4 | 0.00 | 0.00 | 0.00 | 73 |
| 5 | 0.39 | 0.68 | 0.50 | 309 |
| 6 | 0.45 | 0.21 | 0.28 | 101 |
| 7 | 0.44 | 0.38 | 0.41 | 60 |
| 8 | 0.72 | 0.18 | 0.28 | 102 |
| 9 | 0.29 | 0.17 | 0.21 | 143 |
| 10 | 0.33 | 0.01 | 0.02 | 105 |
| 11 | 0.55 | 0.57 | 0.56 | 116 |
| 12 | 0.96 | 0.57 | 0.71 | 261 |
| 13 | 0.00 | 0.00 | 0.00 | 65 |
| 14 | 0.63 | 0.35 | 0.45 | 352 |
| 15 | 0.47 | 0.08 | 0.14 | 108 |
| 16 | 0.53 | 0.69 | 0.60 | 343 |
| 17 | 0.57 | 0.33 | 0.42 | 152 |
| 18 | 0.52 | 0.55 | 0.53 | 325 |
| 19 | 0.46 | 0.30 | 0.36 | 53 |
| 20 | 0.58 | 0.31 | 0.41 | 121 |
| 21 | 0.47 | 0.23 | 0.31 | 105 |
| 22 | 0.25 | 0.88 | 0.39 | 272 |
| | | | | |
| accuracy | | | 0.45 | 4002 |
| macro avg | 0.44 | 0.35 | 0.35 | 4002 |
| weighted avg | 0.49 | 0.45 | 0.42 | 4002 |

**Comments on ResNet50**

- ResNet50 gives relatively good accuracy ,very much higher than VGG but less than Xception
- It gives better than VGG as it is more deep and can extract more abstract features from the images
- Still ,Xception gives more accurate results more than ResNet

**Final result**

From above results of the three models ,Xception is the chosen model due to its performance on the test set and in terms of classification report and test 1 score
Parameters

- Dropout rate =0.7 in the classifier layer
- Learning rate=0.0001
- Model base layers from 30 to 60 are trained

Our final result

| Test accuracy | F1 score | recall | Precision |
|---|---|---|---|
| .6314 | .62 | .63 | .64 |

**Comparing our results to testbenches and state of art solutions**

Our final result

| Test accuracy | F1 score | recall | Precision |
|---|---|---|---|
| .6314 | .62 | .63 | .64 |

- These results are better than all the results on kaggle for dermnet dataset (of course we have the advantage of a little larger training set due to the merging)
- We could not find any paper (after a lot of search ) concerning this particular dataset to compare with.
- Papers discussing similar problems (skin disease detection) achieve Accuracy of around ~80 and 85 % but almost all of them are either on one class detection or fewer classes than ours.

## Ethical issues

The most noticeable ethical issues related to our model are :
- Fairness and biases of our model and implementation
  Because most of the images in our dataset are for white people ,we expect our model to be biased.This problem could be solved by providing data for skin diseases with black people images.
- There is an ethical issue towards other DL specialists. We take it into consideration because we have trained our model on other data to improve its accuracy by increasing its training dataset capacity. To be ethical, we tested our model by using the original testing dataset as it was.

## Discussion

The Xception model is very sophisticated but very complex and deep at the same time.That is why it gives best accuracy out of the three models.

But being a very deep model with a small dataset ,an overfitting problem happens which we tried to reduce by some techniques such as dropout ,L2 regularization and reducing the number of output layers.

Then by noticing the performance ,we finally chose the parameter mentioned above as our final parameters for the Xception model.

For the performance of the model :

Our final result

| Test accuracy | F1 score | recall | Precision |
|---|---|---|---|
| .6314 | .62 | .63 | .64 |

This result is one of the best if compared to testbenches and state of the art.but there are some limitations to the performance of our model such as the size of the dataset is considerably small with 23 classes to predict.

The hard classification task of predicting 23 classes of skin diseases with very similar shapes which tend to confuse the model and resulting in misclassification.

The huge imbalance in the dataset ,which results in degrading the performance of the model.

Things to improve the model performance

- Increasing the size of dataset from a real source of possible (the solution to all problems )
- More Data Augmentation (to increase the size of dataset and
- Increasing the size of dataset by using techniques such as Generative Adversarial Networks (Gans) to generate new images to solve the problem of dataset size and the imbalance.
- We can enhance the image using approaches in computer vision to try to better images.
- Additionally, we can use SIFT to find the image features and then build a bag of words and train the features using deep learning.

## Contribution of each Team member

- We all worked together in developing the notebook
- For the tunning ,each member took some cases to run and reporting the result
- For the results ,each one of us put it in the report
- For comments and discussion ,this done as a teamwork since we needed all our heads together

# References

[1]J. Velasco, "A Smartphone-Based Skin Disease Classification Using MobileNet

CNN," International Journal of Advanced Trends in Computer Science and Engineering,

pp. 2632–2637, Oct. 2019, doi: 10.30534/ijatcse/2019/116852019

[2]P. N. Srinivasu, J. G. SivaSai, M. F. Ijaz, A. K. Bhoi, W. Kim, and J. J. Kang,

"Classification of Skin Disease Using Deep Learning Neural Networks with MobileNet

V2 and LSTM," Sensors, vol. 21, no. 8, p. 2852, Apr. 2021, doi: 10.3390/s21082852.