



## ABSTRACT

Skin cancer is a prevalent and potentially life-threatening disease with increasing incidence worldwide. Early detection plays a crucial role in successful treatment and improved prognosis for patients. This abstract focuses on the development of an advanced skin cancer detection system leveraging cutting-edge technology, such as artificial intelligence and image processing. The proposed system aims to enhance the accuracy and efficiency of skin cancer diagnosis by analyzing dermoscopic images. Utilizing a deep learning approach, the model is trained on a diverse dataset to recognize various skin cancer types, including melanoma, basal cell carcinoma, and squamous cell carcinoma. The project follows a systematic approach, beginning with the collection of a diverse dataset of dermoscopic images representing different skin cancer types. Data preprocessing involves image normalization, enhancement, and segmentation to extract relevant features. Feature extraction is a crucial step to capture essential characteristics of skin lesions, enabling effective differentiation between benign and malignant cases. A comprehensive evaluation of machine learning algorithms is conducted to determine the optimal model for skin cancer classification. Commonly employed algorithms, including Support Vector Machine (SVM), Random Forest, and Convolutional Neural Networks (CNN), are implemented and fine-tuned. The model's performance is assessed using metrics such as accuracy, sensitivity, specificity, and area under the receiver operating characteristic curve (AUC-ROC). The significance of a project focused on skin cancer detection using machine learning is paramount in the realm of healthcare and medical research. Skin cancer, a prevalent and potentially life-threatening disease, underscores the urgent need for advanced diagnostic tools. By leveraging machine learning algorithms to analyze dermoscopic images, this project aims to revolutionize early detection methodologies. Early diagnosis is critical, as it not only significantly improves patient outcomes and reduces mortality rates but also has the potential to alleviate the burden on healthcare systems by enabling more cost-effective interventions.

**Keywords:** Melanoma, Squamous Cell Carcinoma, Convolutional Neural Networks, Support Vector Machine, Accuracy, Sensitivity, Specificity, Interpretability, Medical imaging, Dermoscopic Image.

# LIST OF FIGURES

4.1	<b>General Architecture.</b>	9
4.2	<b>Data Flow Diagram</b>	10
4.3	<b>Use Case Diagram</b>	11
4.4	<b>Class Diagram</b>	12
4.5	<b>Sequence Diagram</b>	13
4.6	<b>Collaboration Diagram</b>	14
4.7	<b>Activity Diagram</b>	15
4.8	<b>Collection of Data</b>	19
4.9	<b>Processing of Data</b>	20
4.10	<b>Decision Tree Algorithm</b>	21
4.11	<b>Logistic Regression</b>	22
4.12	<b>Random Forest Classifier</b>	23
4.13	<b>K-Nearest Neighbours</b>	24
4.14	<b>Google Colab Environment</b>	25
4.15	<b>Loading Data Set</b>	25
4.16	<b>Decision Tree</b>	26
4.17	<b>Logistic Regression</b>	27
4.18	<b>Random Forest Classifier</b>	27
4.19	<b>K-Nearest Neighbour</b>	28
5.1	<b>Permissions to Dataset</b>	29
5.2	<b>Testing ML Models</b>	33
5.3	<b>Test Image</b>	35
6.1	<b>Processed Skin Tissue</b>	40

# **LIST OF ACRONYMS AND ABBREVIATIONS**

CNN	Convolutional Neural Networks
DI	Dermatoscopic Images
ELM	Extreme Learning Machine
IDE	Integrated Development Environment
ML	Machine Learning
PC	Personal Computer
RF	Random Forest
SVM	Support Vector Machine

# TABLE OF CONTENTS

	Page.No
<b>ABSTRACT</b>	<b>2</b>
<b>LIST OF FIGURES</b>	<b>3</b>
<b>LIST OF ACRONYMS AND ABBREVIATIONS</b>	<b>4</b>
<b>1 INTRODUCTION</b>	<b>1</b>
1.1 Introduction . . . . .	1
1.2 Aim of the Project . . . . .	1
1.3 Project Domain . . . . .	2
1.4 Scope of the Project . . . . .	2
<b>2 LITERATURE REVIEW</b>	<b>4</b>
<b>3 PROJECT DESCRIPTION</b>	<b>6</b>
3.1 Existing System . . . . .	6
3.2 Proposed System . . . . .	6
3.3 Feasibility Study . . . . .	7
3.3.1 Economic Feasibility . . . . .	7
3.3.2 Technical Feasibility . . . . .	7
3.3.3 Social Feasibility . . . . .	7
3.4 System Specification . . . . .	8
3.4.1 Hardware Specification . . . . .	8
3.4.2 Software Specification . . . . .	8
3.4.3 Standards and Policies . . . . .	8
<b>4 METHODOLOGY</b>	<b>9</b>
4.1 General Architecture . . . . .	9
4.2 Design Phase . . . . .	10
4.2.1 Data Flow Diagram . . . . .	10
4.2.2 Use Case Diagram . . . . .	11
4.2.3 Class Diagram . . . . .	12
4.2.4 Sequence Diagram . . . . .	13
4.2.5 Collaboration Diagram . . . . .	14
4.2.6 Activity Diagram . . . . .	15

4.3	Algorithm & Pseudo Code . . . . .	16
4.3.1	Algorithm . . . . .	16
4.3.2	Pseudo Code . . . . .	16
4.4	Module Description . . . . .	19
4.4.1	Collection Of Data: . . . . .	19
4.4.2	Processing Of Data: . . . . .	20
4.4.3	Implementing Algorithms: . . . . .	21
4.5	Steps to execute/run/implement the project . . . . .	25
4.5.1	Run Python : . . . . .	25
4.5.2	Loading Data Set: . . . . .	25
4.5.3	Decision Tree: . . . . .	26
<b>5</b>	<b>IMPLEMENTATION AND TESTING</b>	<b>29</b>
5.1	Input and Output . . . . .	29
5.1.1	Permissions to Dataset . . . . .	29
5.1.2	Processing of the Dermal . . . . .	30
5.2	Testing . . . . .	30
5.3	Types of Testing . . . . .	30
5.3.1	Unit Testing . . . . .	30
5.3.2	Integration Testing . . . . .	32
5.3.3	System Testing . . . . .	33
5.3.4	Test Result . . . . .	35
<b>6</b>	<b>RESULTS AND DISCUSSIONS</b>	<b>36</b>
6.1	Efficiency of the Proposed System . . . . .	36
6.2	Comparison of Existing and Proposed System . . . . .	36
6.3	Sample Code . . . . .	37
<b>7</b>	<b>CONCLUSION AND FUTURE ENHANCEMENTS</b>	<b>41</b>
7.1	Conclusion . . . . .	41
7.2	Future Enhancements . . . . .	42
<b>8</b>	<b>SOURCE CODE &amp; POSTER PRESENTATION</b>	<b>43</b>
8.1	Source Code . . . . .	43
	<b>References</b>	<b>46</b>

# Chapter 1

## INTRODUCTION

### 1.1 Introduction

Skin cancer represents a substantial global health challenge, with its incidence on the rise. Early detection is paramount for effective intervention and improved prognosis. In recent years, the intersection of medical imaging and machine learning has emerged as a promising avenue for enhancing diagnostic capabilities. This project delves into the development of a skin cancer detection system using machine learning, specifically tailored to analyze dermatoscopic images.

Leveraging advanced algorithms, this system aims to discriminate between various skin cancer types, such as melanoma, basal cell carcinoma, and squamous cell carcinoma, with a focus on achieving high accuracy. The integration of machine learning in dermatology holds the potential to not only expedite the diagnostic process but also to contribute to a paradigm shift in healthcare, enabling timely and precise interventions that can significantly impact patient outcomes.

The primary technical goal of skin cancer detection using machine learning is to develop a robust and accurate algorithm that can analyze dermatoscopic images to automatically identify and classify different types of skin lesions, including benign and malignant cases.

The machine learning model aims to learn complex patterns and features from a diverse dataset, enabling it to make precise predictions about the nature of skin abnormalities. This involves implementing advanced image processing techniques, feature extraction methods, and classification algorithms to achieve high sensitivity and specificity in detecting various forms of skin cancer, such as melanoma, basal cell carcinoma, and squamous cell carcinoma.

### 1.2 Aim of the Project

The aim of our project, the overarching aim of the project on skin cancer detection using machine learning is to pioneer a robust and accurate diagnostic tool for the early and precise identification of skin lesions. Leveraging the power of advanced

machine learning algorithms, the project seeks to develop a system capable of analyzing dermoscopic images and autonomously classifying various types of skin cancer, including melanoma, basal cell carcinoma, and squamous cell carcinoma. The central focus is on enhancing diagnostic accuracy, minimizing false positives and negatives, and ultimately improving patient outcomes through timely interventions. Beyond accuracy, the project aims to create a seamlessly integrated solution, ensuring practicality and user-friendliness for healthcare professionals.

The primary aim of the project on skin cancer detection using machine learning is to harness advanced computational techniques to create a robust and efficient system for the early and accurate identification of skin lesions. Grounded in the urgency of early intervention, the project seeks to leverage machine learning algorithms to analyze dermoscopic images, enabling precise categorization of skin abnormalities, including various types of skin cancer like melanoma, basal cell carcinoma, and squamous cell carcinoma. The overarching goal is to enhance the diagnostic process, providing healthcare professionals with a reliable tool that minimizes false positives and negatives, ultimately improving patient outcomes.

### **1.3 Project Domain**

The project on skin cancer detection using machine learning operates within the domain of medical imaging and healthcare technology. Specifically, it addresses the critical intersection of machine learning and dermatology. This domain involves the utilization of advanced computational techniques to analyze dermoscopic images, with the primary objective of developing a sophisticated diagnostic tool for the early detection and classification of skin cancer. Grounded in the urgency of timely intervention, the project navigates the intricacies of medical image analysis, feature extraction, and classification algorithms to create a system capable of autonomously identifying various types of skin lesions, ranging from benign to malignant.

### **1.4 Scope of the Project**

The scope of the project on skin cancer detection using machine learning is vast and encompasses various crucial dimensions in the intersection of healthcare and technology. The primary focus is on developing an advanced system that utilizes machine learning algorithms to analyze dermoscopic images for the accurate and early detection of skin cancer. This involves delving into the intricacies of medical image analysis, with an emphasis on feature extraction, pattern recognition, and classification of diverse skin lesions. The project aims to not only distinguish between



benign and malignant cases but also to categorize specific skin cancer subtypes, such as melanoma, basal cell carcinoma, and squamous cell Carcinoma.

The essence of the project lies in revolutionizing the landscape of skin cancer detection through the integration of machine learning. At its core, the project seeks to develop a sophisticated and accurate diagnostic tool capable of autonomously analyzing dermoscopic images to identify and classify various types of skin lesions. The heart of the project beats with a commitment to early detection, aiming to enhance patient outcomes by identifying skin cancer at its incipient stages.

This involves navigating the intricacies of medical image analysis, delving into advanced machine learning algorithms, and optimizing their performance for robust and efficient detection.

Proactive maintenance strategies are pivotal for sustaining the effectiveness and reliability of a skin cancer detection system based on machine learning. Regular model retraining stands as a cornerstone, ensuring the algorithm remains updated with contemporary data, adapting to emerging patterns in dermoscopic images.

Continuous data monitoring is imperative, allowing for the early identification of any shifts or biases in input data, which is essential for maintaining the model's robustness. Dynamic feature engineering ensures that feature extraction methods evolve with the latest advancements in image processing technology, enabling the system to capture relevant diagnostic information effectively.

Enhancements in interpretability contribute to a deeper understanding of the model's decisions, fostering trust and consistent utilization by healthcare professionals. Scalability considerations, user training, and compliance with regulatory standards round out the proactive measures, ensuring that the skin cancer detection system remains reliable, secure, and aligned with evolving clinical and technological landscapes.

## Chapter 2

# LITERATURE REVIEW

[1] Jianhua Zhao et al, proposed the Mobile Data Management. This paper describes an advanced seminar presented at the 16th IEEE International Conference on Mobile Data Management. The advanced seminar presents the state-of-the-art in the Field of Skin Cancer Detection, which is fast emerging as a disruptive technology for years to come. The seminar focuses on Internet-scale sensor information management, related mobile analytics open source ML technologies, and emerging standards.

[2] Haishan Zeng et al, addressed the Skin Cancer Detection. This project introduces a smart waste bin employing machine learning to predict Skin Cancer detection. The system alerts authorities when bins reach critical levels, optimizing types of Skin Cancer Sample collection routes. Continuous data analysis, presented through cloud-based graphs, enhances monitoring. Automatic Skin detection notify authorities, saving time and resources, minimizing Skin Cancer, and preventing the spread of diseases associated with Different type of skin Disease.

[3]David I. McWilliams et al, explained the Real life Smart Skin Cancer Detection Management System. This paper explores automated machine learning in a real-life Smart Skin Cancer Detection Management system, focusing on binary classification of Skin problem using Machine Learning. Various methods, including modified manual models and conventional algorithms, are compared. The best solution, a Random Forest classifier, significantly enhances accuracy and recall,improving problem face by the skin cancer detection.

[4] Wenbo Wang et al, illustrates the Extreme Learning Machine. This paper addresses the slow learning speed of feedforward neural networks, identifying two key reasons: the use of slow gradient-based learning algorithms and iterative tuning of all network parameters. It introduces a novel learning algorithm, Extreme Learning Machine (ELM), for single-hidden layer feedforward neural networks. ELM randomly selects input weights and analytically determines output weights, theoretically offering optimal generalization performance with exceptionally fast learning speeds. Experimental results demonstrate superior generalization and faster learning compared to traditional algorithms across various real-world problems.

[5] Agung W. Setiawan et al, illustrates the Transporting and Causing Property Damage. Hong Kong faces challenges during the rainy season, impacting transportation and causing property damage. Limited research on drainage systems prompts the

exploration of IoT for smart city development. A prototype system integrates IoT, collecting data for Artificial Neural Network training. Proposed predictive maintenance solutions for skin cancer detection. Well-trained algorithms effectively predict drainage situations, offering reliability and societal benefits for Hong Kong's drainage services.

[6] Pratyaksh P et al, addressed the Skin Problem Levels in India. In response to escalating Skin problem levels in India, an ML-based smart monitoring system is proposed. A deep learning model predicts garbage levels with 80.33 percentage accuracy, facilitating efficient waste management. Results confirm precise prognosis, advocating for the integration of ML and deep learning to revolutionize technology in Cancer management overflow.

[7] Qinghua Wang et al, proposed the Automating Data Collection and Monitoring. Ongoing digital transformation in skin management involves automating data collection and monitoring. A smart skin cancer system in southern Sweden, with deployed sensors, enables real-time data analysis. Preliminary results indicate the development of statistical models for hospital, society facilitating data forecasting and anomaly detection.

[8] Z. Zhang et al, illustrates the Importance of Accurate Data for Situational Awareness. This study addresses inflow and infiltration (I/I) in skin management systems, emphasizing the importance of accurate data for situational awareness. Employing AI models (ANFIS and MLPNN) and SCADA system data, the research reveals both subcatchments exhibit I/I. ANFIS outperforms MLPNN in modeling I/I situations, aiding spatial decision-making for cancer stage system maintenance.

## Chapter 3

# PROJECT DESCRIPTION

### 3.1 Existing System

Deployment of skin cancer detection using machine learning have showcased advancements in the field of medical image analysis. These systems leverage sophisticated algorithms to analyze dermatoscopic images for the early identification and classification of skin lesions. Notable approaches include the application of deep learning techniques, such as Convolutional Neural Networks (CNNs), for feature extraction and pattern recognition.

Many existing systems focus on multi-class classification, distinguishing between benign and malignant lesions and categorizing specific subtypes, including melanoma, basal cell carcinoma, and squamous cell carcinoma. The integration of large and diverse datasets is a common practice to enhance the generalizability of these systems across different skin types and clinical scenarios.

Interpretability features have gained attention to provide insights into the decision-making process of machine learning models, fostering trust among healthcare professionals. These existing systems have shown promising results, ongoing research strives to address challenges such as limited interpretability, diverse dataset representation, and seamless integration into clinical workflows, aiming to further improve the accuracy and practicality of skin cancer detection using machine learning. It's essential to consult the latest literature for the most recent developments and advancements in this rapidly evolving field..

### 3.2 Proposed System

The proposed system for skin cancer detection using machine learning aims to usher in a new era of precision and efficiency in dermatological diagnostics. Employing cutting-edge deep learning architectures, such as Convolutional Neural Networks (CNNs), the system focuses on extracting intricate features from dermatoscopic images, enhancing its ability to discriminate between benign and malignant lesions.

The algorithm for skin cancer detection using machine learning is a systematic process designed to effectively analyze Dermatoscopic Images and distinguish between benign and malignant skin lesions. It commences with the collection of a di-

verse and comprehensive dataset, containing images representing various skin types and lesion characteristics. Following data preprocessing steps, such as normalization and resizing, the dataset is divided into training and testing sets.

The model is trained on the labeled dataset, where it learns to differentiate between benign and malignant lesions. Subsequently, the algorithm undergoes validation to fine-tune parameters and prevent overfitting, followed by testing on an independent dataset to evaluate its real-world performance. Interpretability features, such as Grad-CAM, may be integrated to provide visual explanations of the model's decisions.

Post-processing steps, including threshold adjustments, refine the model's outputs. The algorithm culminates in the integration of the developed model into the clinical workflow, ensuring that healthcare professionals can seamlessly leverage its capabilities for accurate and timely skin cancer diagnosis. Continuous monitoring and updating mechanisms guarantee the adaptability of the algorithm to evolving skin cancer patterns, reinforcing its reliability in clinical practice.

### **3.3 Feasibility Study**

#### **3.3.1 Economic Feasibility**

The project costs are optimal because the requirements are a PC, data about the previous overflow (we can either gather the data and create a dataset but as it requires more data and a dataset is already used), Python IDE, and MSoffice.

#### **3.3.2 Technical Feasibility**

The minimum requirements are Windows 7/8/10/11, storage space, and a processor of at least dual core, so technically the project runs fine without any heavy load process.

#### **3.3.3 Social Feasibility**

Since the project is both economically and technically feasible, and the user has easy access to the system, the software is very well socially feasible.

### **3.4 System Specification**

#### **3.4.1 Hardware Specification**

A computer with

- 1 . Processor : Pentium IV or higher
- 2 . RAM : 2GB
- 3 . Space on Disk : minimum 512 MB.

#### **3.4.2 Software Specification**

1. For developing the application the python using Google Colab and Arduino IDE are used.
2. Operation systems Supported : Windows 7, Windows 8, Windows 10, Windows 11.

#### **3.4.3 Standards and Policies**

1. Python standard Used: Google Colab.
2. CNN Algorithm.

##### **Google Colab**

Google Colab is a type of command line interface which explicitly deals with the ML( MachineLearning) modules.And navigator is available in all the Windows,Linux and MacOS.The Google Colab has many number of IDE's which make the coding easier. The UI can also be implemented in python.

##### **Standard Used: ISO/IEC 27001**

**CNN Algorithm** A Convolutional Neural Network (CNN) is a type of artificial neural network designed specifically for processing and analyzing visual data. It has proven highly effective in tasks such as image classification, object detection, and image generation.

# Chapter 4

## METHODOLOGY

### 4.1 General Architecture

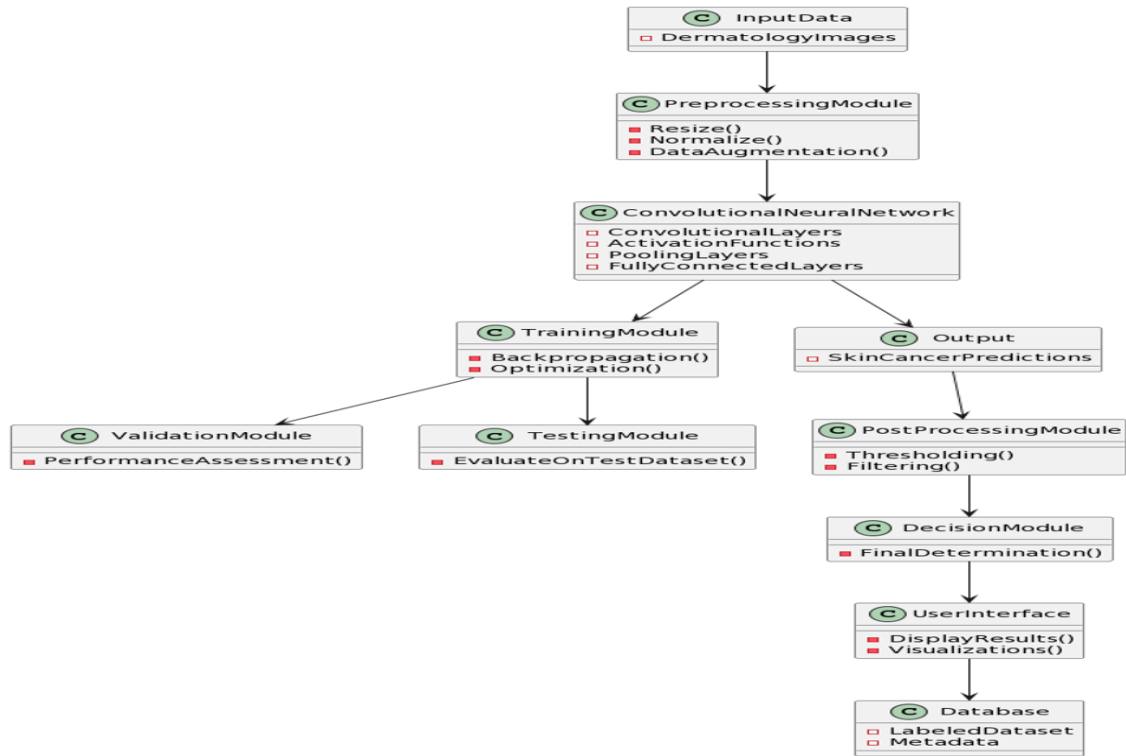


Figure 4.1: General Architecture.

The figure 4.1 shows system architecture is designed with a comprehensive approach for skin cancer detection. After that, we will clean the data to obtain the required information. Next, we will perform data preprocessing, which involves cleaning, transforming, and organizing the raw data to improve its quality. We will then visualize the collected data using graphs, which makes it easy to interpret and reveal patterns for efficient decision-making. We will also quantify the degree of linear relationship between two variables in the data using the correlation coefficient, which ranges from -1 to 1.

## 4.2 Design Phase

### 4.2.1 Data Flow Diagram

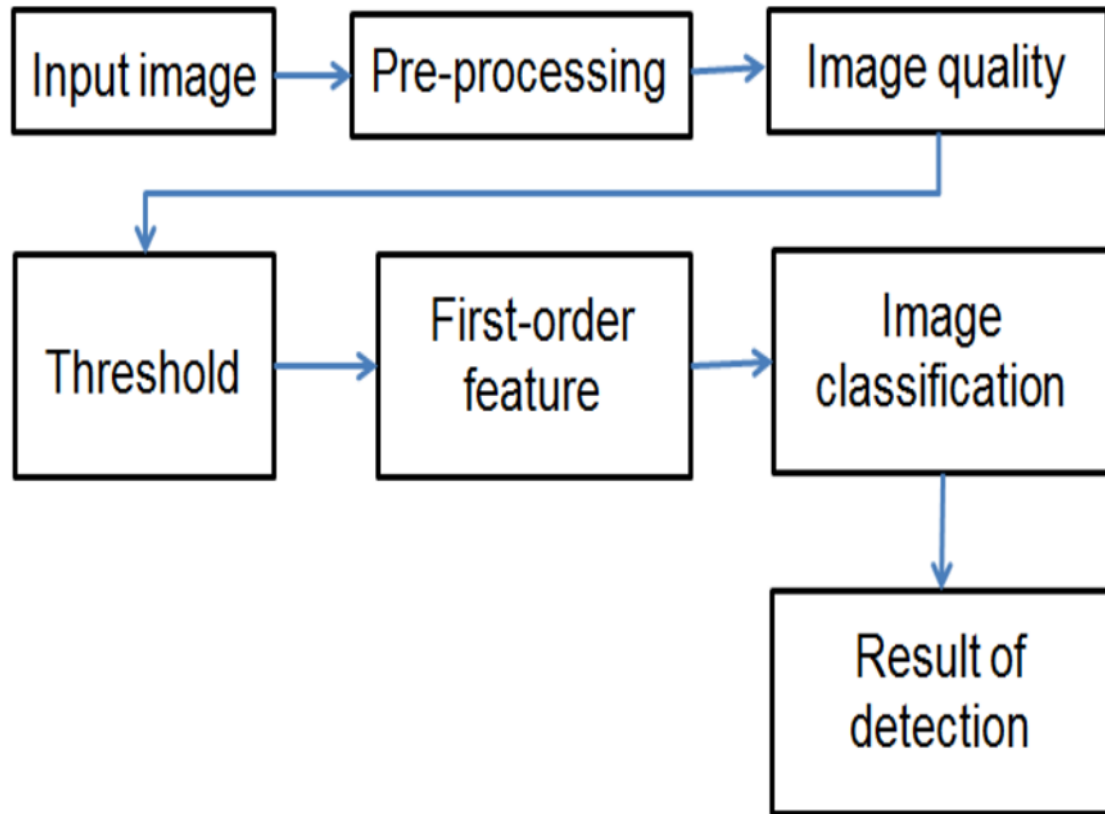


Figure 4.2: **Data Flow Diagram**

The figure 4.2 shows the data flow diagram to understand systematic progression of information in the skin cancer deection. The data collected from previous tests go through cleaning and preprocessing and saved in .csv format. ML algorithms are applied to train the data and test it to check if the trained dataset works properly for future impacts and gets used. Then they are applied in the real-time application.



#### 4.2.2 Use Case Diagram

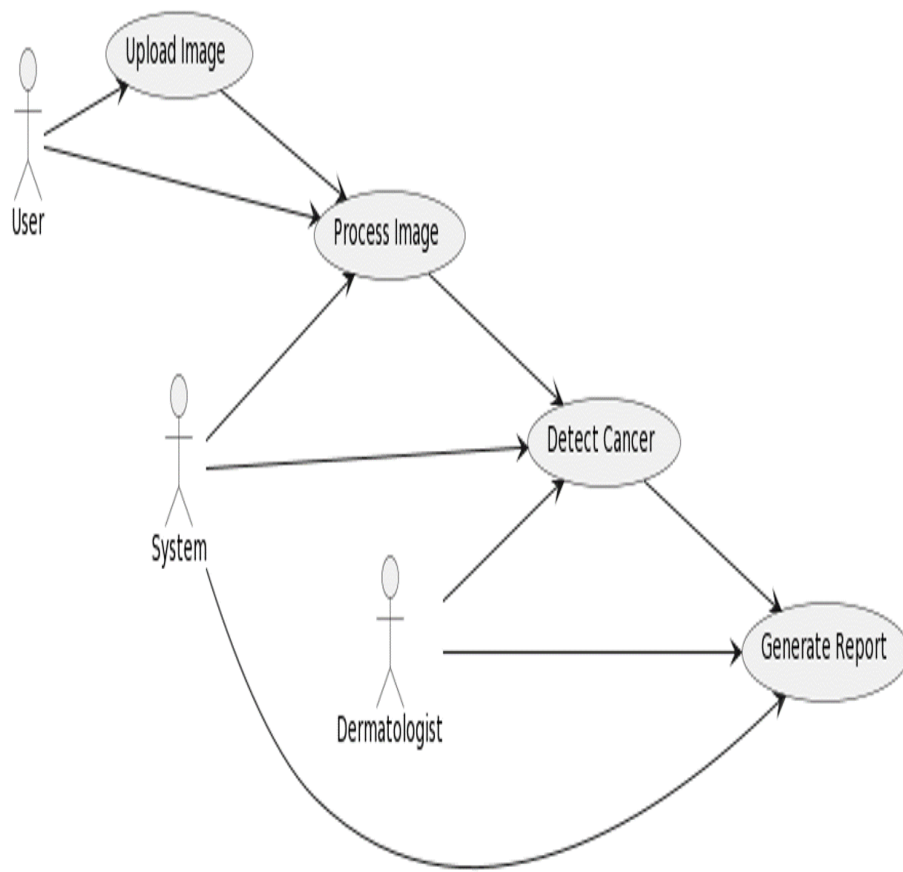


Figure 4.3: Use Case Diagram

The figure 4.3 shows the use case diagram to illustrate the interactions between them. The real-time monitor installed in the detector App system will constantly monitor the flow to detect any overflows. If any Memory of Data overflows are detected, the cloud-based monitoring system will notify the user with the details of the issue. The user will then assign an employee working nearby to clear the overflow and unblock the system.

### 4.2.3 Class Diagram

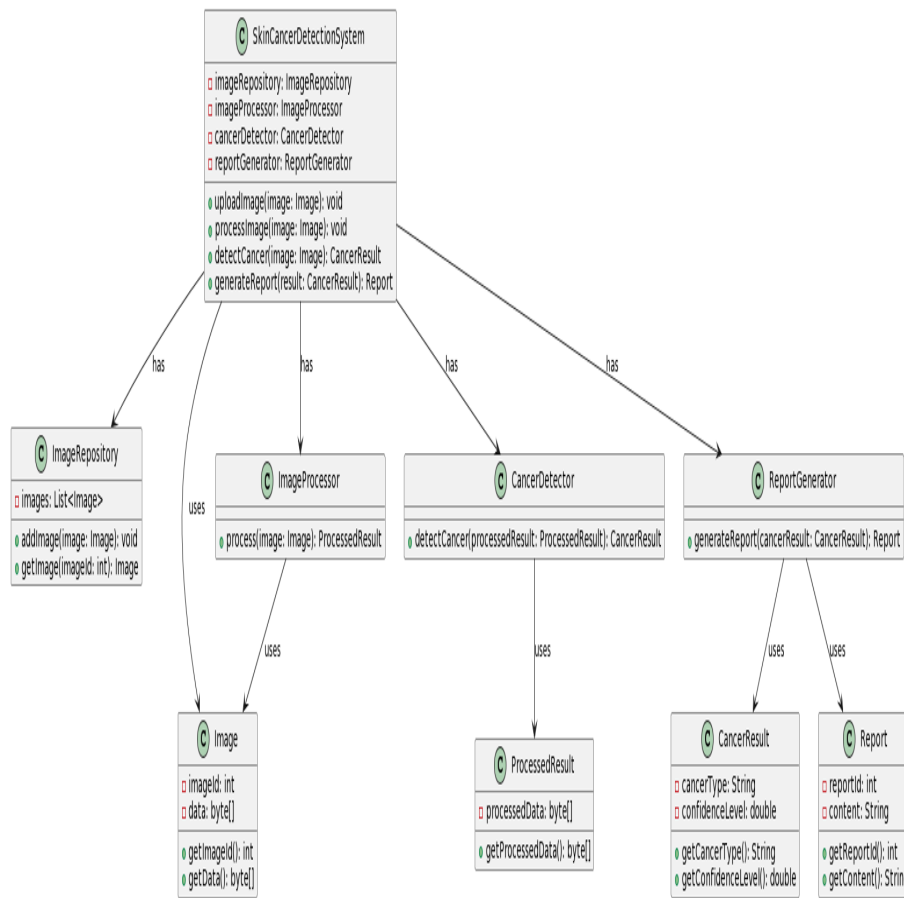


Figure 4.4: Class Diagram

The figure 4.4 shows the class diagram which outlines the fundamental components of the skin cancer detection. The dataset contains information about the man-hole and the sensor that are stored in a database. This database is monitored by the cloud, which sends notifications to the User in the event of a Detection. The User will then assign an employee to inspect the area and clear the Memory.

#### 4.2.4 Sequence Diagram

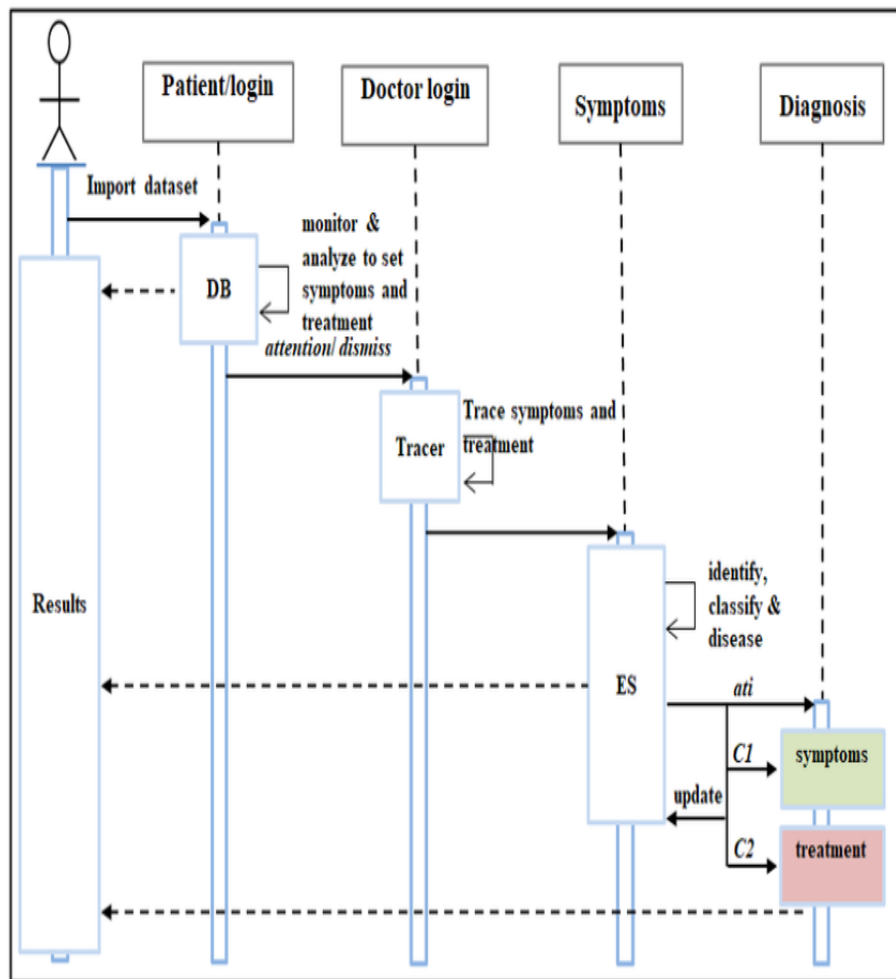


Figure 4.5: Sequence Diagram

The figure 4.5 shows the sequence diagram which illustrates the dynamic interactions within the system. Sequence diagram represents the passage of messages.

**Cloud:** is mode of a computer data storage in which digital data is stored on servers in off site locations. Helps in various data base operations. It monitors the data periodically and alert the database host provider in cases of any anomalies detected.

**ML algorithms:** Machine Learning algorithms are used to calculate the accuracy of readings taken out. Machine Learning Library contain many algorithms that are used to derive the best accuracy from the read data.

**Accuracy:** states the Accuracy of reading, that will be calculated by using the appropriate machine learning algorithms such as Random Forest Classifier, Logistic regression, etc.

**User:** Cloud will present the calculated Accuracy to the authorized members of the User. User receive alert of the Skin problem.

#### 4.2.5 Collaboration Diagram

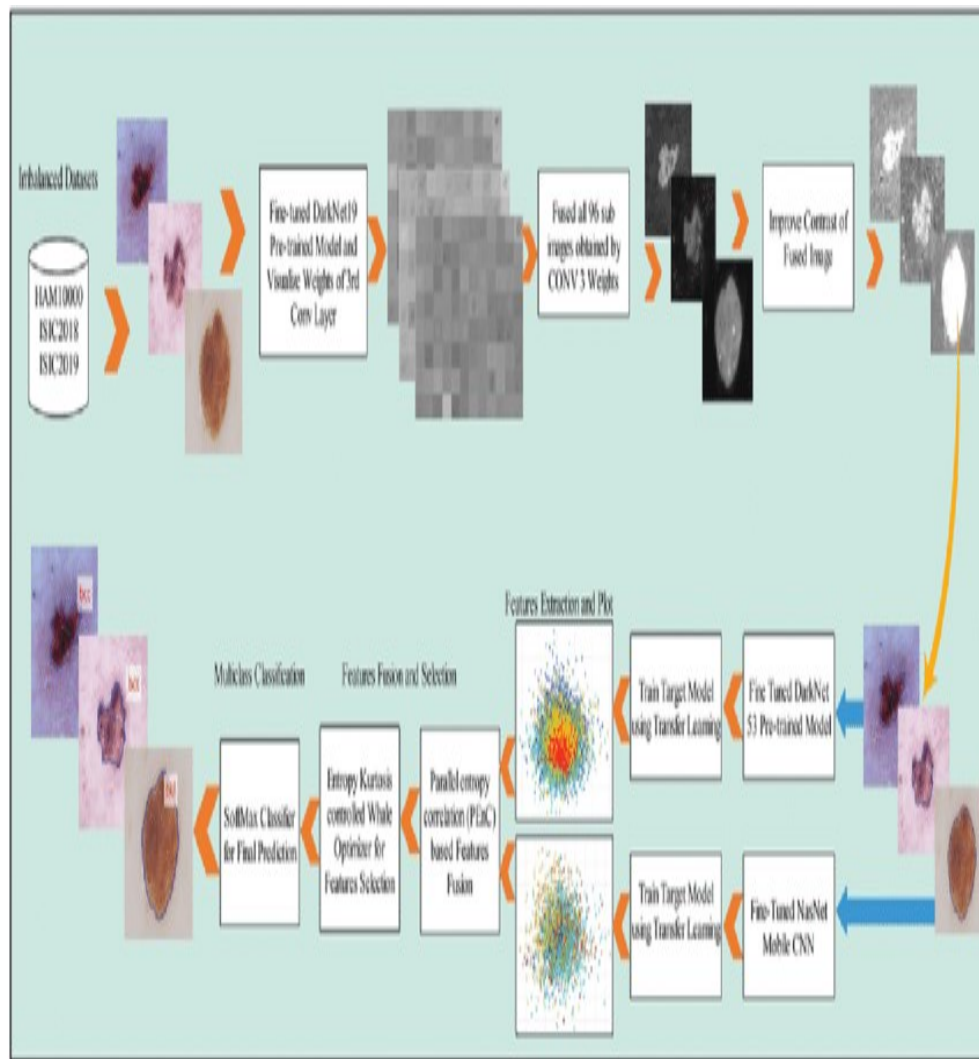


Figure 4.6: Collaboration Diagram

The figure 4.6 shows the collaboration diagram which depicts the structure of the project. The logistic regression algorithm is mainly used for classification tasks where the goal is to predict the probability that an instance belongs to a given class or not. A decision tree is one of the most powerful tools used for both classification and regression tasks. It builds a flowchart-like tree structure where each internal node denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node (terminal node) holds a class label. The K-Nearest Neighbors (KNN) algorithm is a robust and intuitive machine learning method employed to tackle classification and regression problems. By capitalizing on the concept of similarity, KNN predicts the label or value of a new data point by considering its K closest neighbors in the training dataset. The SVM method draws a graph and takes the plane with the most area.

#### 4.2.6 Activity Diagram

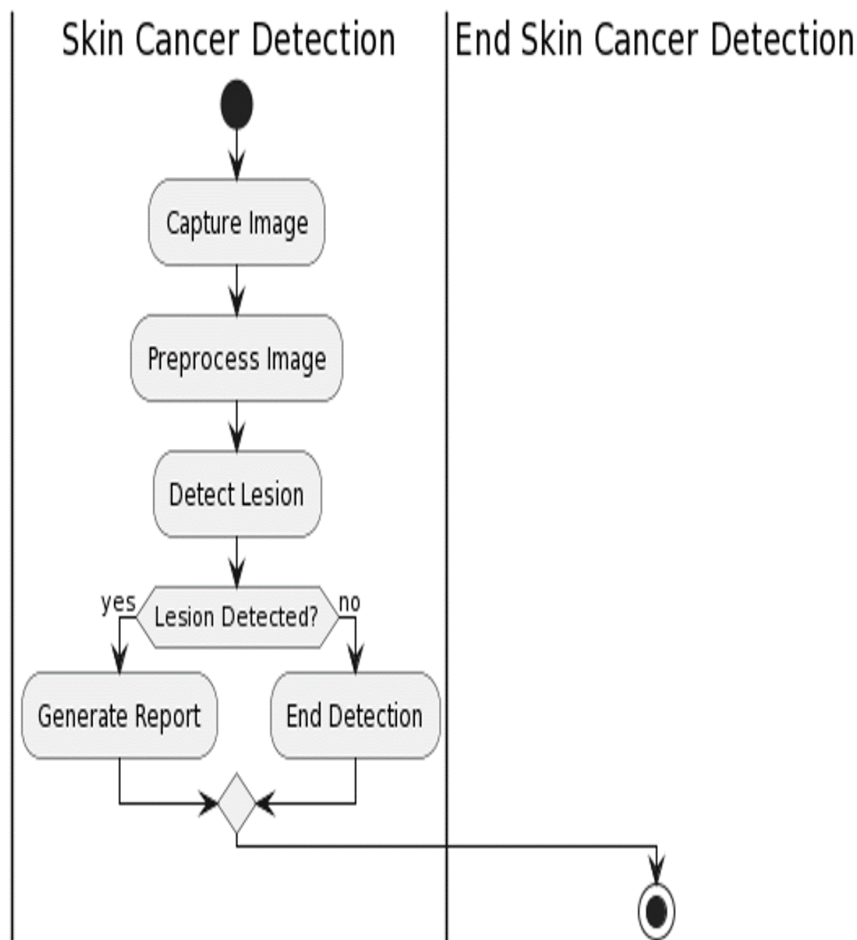


Figure 4.7: Activity Diagram

The figure 4.7 shows the activity diagram which depicts the flow of events of skin cancer detection. First, data is collected and performed preprocessing. The data is trained info is tested and checked whether blocking has occurred or not. If the blocking, the idea is to inform the user so that they will assign to user to skin problem, if not there would be no problem.

## 4.3 Algorithm & Pseudo Code

### 4.3.1 Algorithm

**Step 1:** Start.

**Step 2:** Clean and load the dataset into the code.

**Step 3:** Generate and train the dataset using Machine Learning Algorithms.

**Step 4:** Run the unit code.

**Step 5:** Visualize the efficiency and accuracy of the graph.

**Step 6:** Stop.

### 4.3.2 Pseudo Code

```
1
2 # Import necessary libraries
3 import pandas as pd
4 import numpy as np
5 import seaborn as sns
6 import plotly.express as px
7 import matplotlib.pyplot as plt
8 from sklearn.preprocessing import StandardScaler
9 from sklearn.model_selection import train_test_split
10 from sklearn.metrics import accuracy_score, confusion_matrix
11 from sklearn.linear_model import LogisticRegression
12 from sklearn.tree import DecisionTreeClassifier
13 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
14 from sklearn.neighbors import KNeighborsClassifier
15 from sklearn.svm import SVC
16 from xgboost import XGBClassifier
17
18 # Load the dataset
19 df = pd.read_csv("/content/drive/Colab Notebooks/Minor 1/waterpotability.csv")
20
21 # Data preprocessing and visualization
22 # ...
23
24 # Handle missing values
25 # ...
26
27 # Feature scaling
28 # ...
29
30 # Split the dataset into training and testing sets
31 # ...
32
33 # Model training and evaluation
34
```

```

35 # Logistic Regression
36 model_lr = LogisticRegression()
37 model_lr.fit(x_train, y_train)
38 pred_lr = model_lr.predict(x_test)
39 accuracy_lr = accuracyscore(y_test, pred_lr)
40
41 # Decision Tree
42 model_dt = DecisionTreeClassifier(max_depth=4)
43 model_dt.fit(x_train, y_train)
44 pred_dt = model_dt.predict(x_test)
45 accuracy_dt = accuracy_score(y_test, pred_dt)
46
47 # Random Forest
48 model_rf = RandomForestClassifier()
49 model_rf.fit(x_train, y_train)
50 pred_rf = model_rf.predict(x_test)
51 accuracy_rf = accuracyscore(y_test, pred_rf)
52
53 # K-Nearest Neighbour
54 model_knn = K-NearestClassifier(n_neighbors=11)
55 model_knn.fit(x_train, y_train)
56 pred_knn = model_knn.predict(x_test)
57 accuracy_knn = accuracyscore(y_test, pred_knn)
58
59 # Support Vector Machine (SVM)
60 model_svm = SVC(kernel="rbf")
61 model_svm.fit(x_train, y_train)
62 pred_svm = model_svm.predict(x_test)
63 accuracy_svm = accuracyscore(y_test, pred_svm)
64
65 # AdaBoost
66 model_ada = AdaBoostClassifier(n_estimators=200, learning_rate=0.03)
67 model_ada.fit(x_train, y_train)
68 pred_ada = model_ada.predict(x_test)
69 accuracy_ada = accuracyscore(y_test, pred_ada)
70
71 # XGBoost
72 model_xgb = XGBClassifier(n_estimators=100, learning_rate=0.04)
73 model_xgb.fit(x_train, y_train)
74 pred_xgb = model_xgb.predict(x_test)
75 accuracy_xgb = accuracyscore(y_test, pred_xgb)
76
77 # Create a dataframe to store model accuracies
78 models = pd.DataFrame({
79     "Model": ["Logistic Regression", "Decision Tree", "Random Forest",
80             "KNN", "SVM", "AdaBoost", "XGBoost"],
81     "Accuracy": [accuracy_lr, accuracy_dt, accuracy_rf,
82                 accuracy_knn, accuracy_svm, accuracy_ada, accuracy_xgb]
83 })
84 # Logistic Regression

```

```

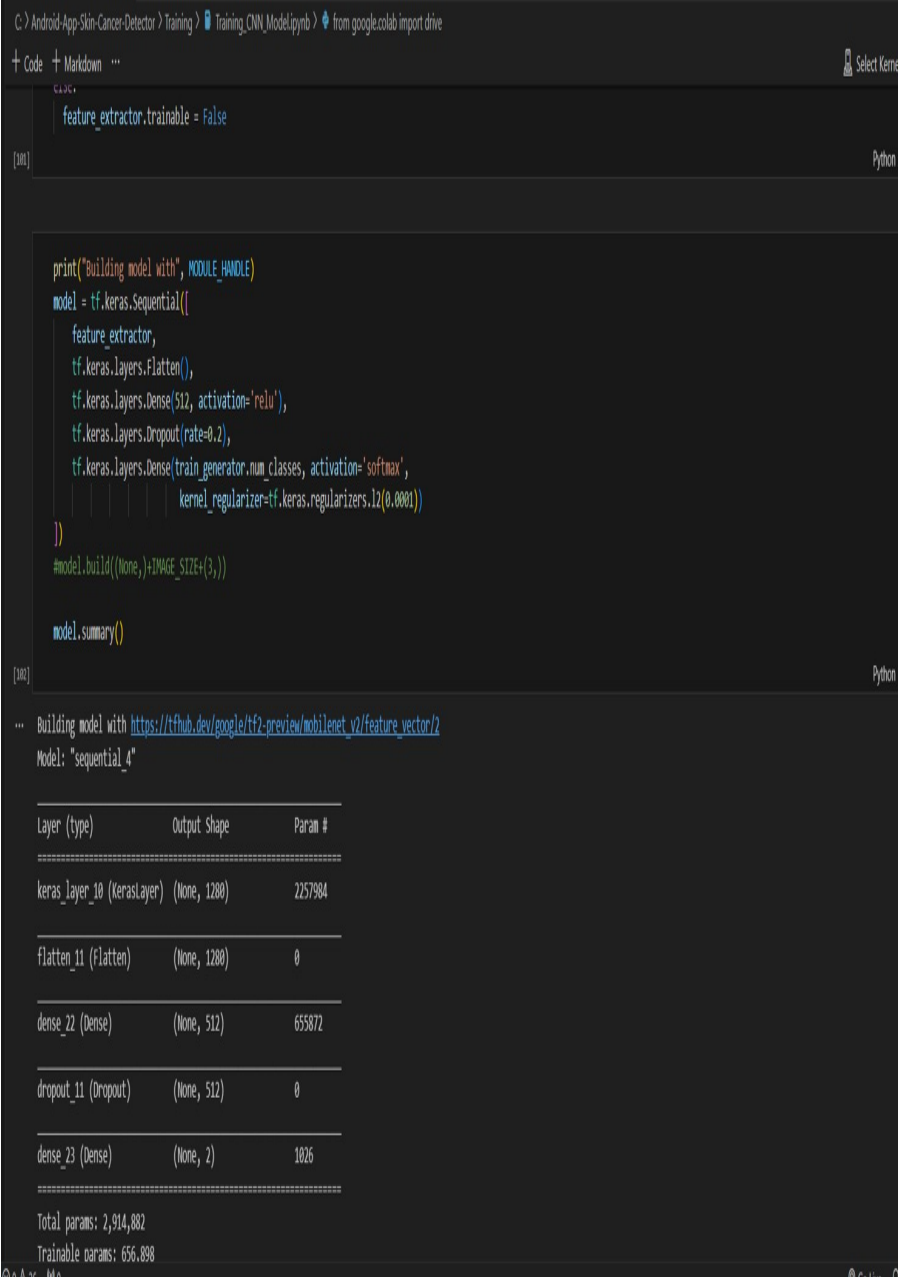
85 model_lr = LogisticRegression()
86 model_lr.fit(x_train, y_train)
87 pred_lr = model_lr.predict(x_test)
88 accuracy_lr = accuracy_score(y_test, pred_lr)
89
90 # Decision Tree
91 model_dt = DecisionTreeClassifier(max_depth=4)
92 model_dt.fit(x_train, y_train)
93 pred_dt = model_dt.predict(x_test)
94 accuracy_dt = accuracy_score(y_test, pred_dt)
95
96 # Random Forest
97 model_rf = RandomForestClassifier()
98 model_rf.fit(x_train, y_train)
99 pred_rf = model_rf.predict(x_test)
100 accuracy_rf = accuracy_score(y_test, pred_rf)
101
102 # K-Nearest Neighbour
103 model_knn = KNeighborsClassifier(n_neighbors=11)
104 model_knn.fit(x_train, y_train)
105 pred_knn = model_knn.predict(x_test)
106 accuracy_knn = accuracy_score(y_test, pred_knn)
107
108 # Visualize model accuracies
109 # ...
110
111 # Display sorted model accuracies
112 # ...

```



## 4.4 Module Description

### 4.4.1 Collection Of Data:



```
C:\Android-App-Skin-Cancer-Detector> Training > Training_CNN_Model.ipynb > from google.colab import drive
+ Code + Markdown ... Select Kernel

feature_extractor.trainable = False

[101] Python

print("Building model with", MODULE_HANDLE)
model = tf.keras.Sequential([
    feature_extractor,
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(512, activation='relu'),
    tf.keras.layers.Dropout(rate=0.2),
    tf.keras.layers.Dense(train_generator.num_classes, activation='softmax',
                           kernel_regularizer=tf.keras.regularizers.L2(0.0001))
])
#model.build((None,)+IMAGE_SIZE+(3,))

model.summary()

[102] Python

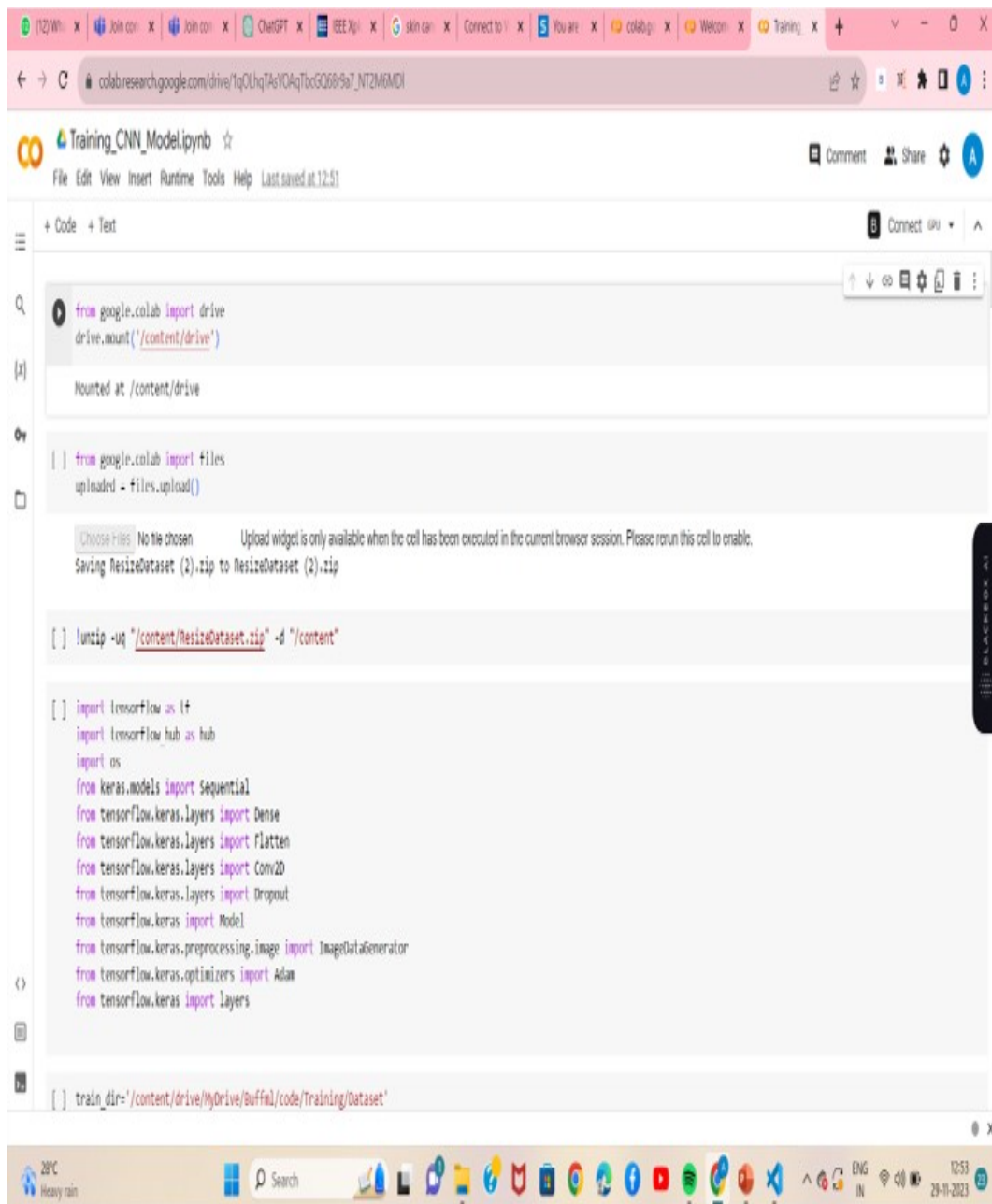
... Building model with https://tfhub.dev/google/tf2-preview/mobilenet_v2/feature_vector/2
Model: "sequential_4"

Layer (type)                 Output Shape          Param #
=====
keras_layer_10 (KerasLayer)  (None, 1280)         2257984
flatten_11 (Flatten)        (None, 1280)          0
dense_22 (Dense)             (None, 512)          655872
dropout_11 (Dropout)         (None, 512)           0
dense_23 (Dense)             (None, 2)             1026
=====
Total params: 2,914,882
Trainable params: 656,898
```

Figure 4.8: Collection of Data

The figure 4.8 shows the Collection of Data .Data collection is the systematic process of gathering and recording information from various sources, such as surveys, observations, or sensing algorithm. It forms the foundation for meaningful analysis and insights in research, decision-making, and problem-solving.

#### 4.4.2 Processing Of Data:



```
from google.colab import drive
drive.mount('/content/drive')

Mounted at /content/drive

from google.colab import files
uploaded = files.upload()

Choose files | No file chosen | Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.
Saving ResizedDataset (2).zip to ResizedDataset (2).zip

[ ] !unzip -uq "/content/ResizedDataset.zip" -d "/content"

[ ] import tensorflow as tf
import tensorflow.hub as hub
import os
from keras.models import Sequential
from tensorflow.keras.layers import Dense
from tensorflow.keras.layers import Flatten
from tensorflow.keras.layers import Conv2D
from tensorflow.keras.layers import Dropout
from tensorflow.keras import Model
from tensorflow.keras.preprocessing.image import ImageDataGenerator
from tensorflow.keras.optimizers import Adam
from tensorflow.keras import layers

[ ] train_dir = '/content/drive/MyDrive/Buffal/code/Training/Dataset'
```

Figure 4.9: Processing of Data

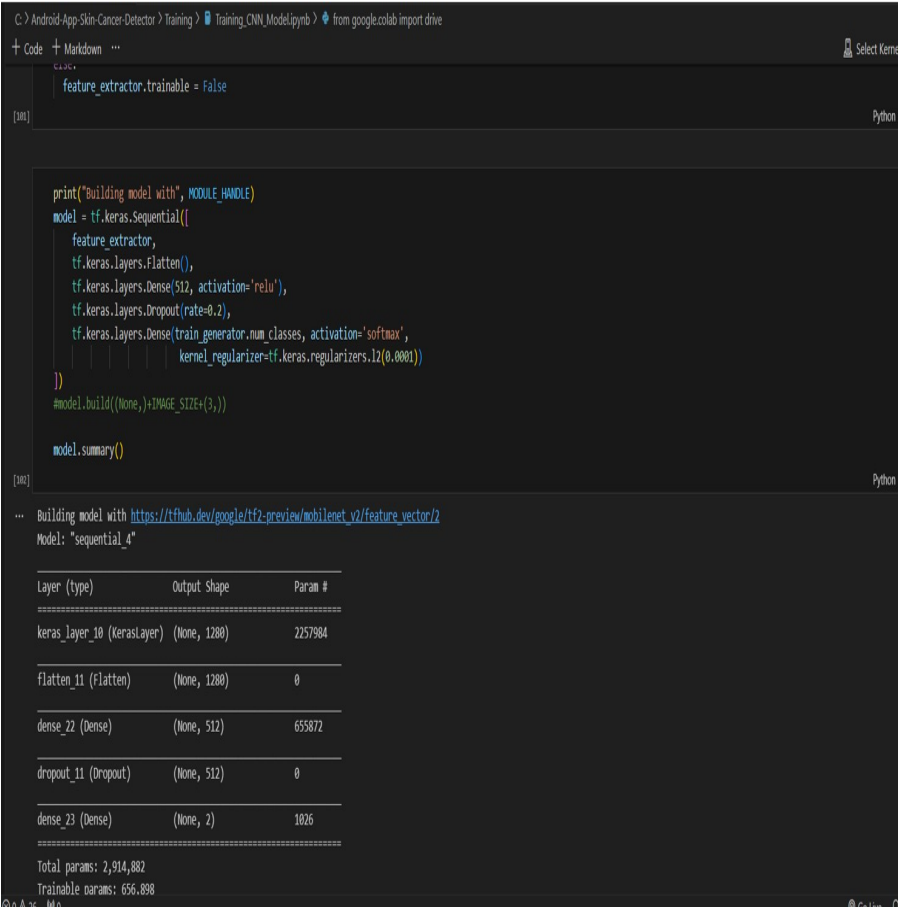
The figure 4.9 shows the Processing of Data .Data preprocessing involves cleaning, transforming, and organizing raw data to enhance its quality and usability. Common tasks include handling missing values, scaling features, and encoding categorical variables, ensuring the data is suitable for analysis or machine learning applications.

## Decision Tree:



21

## Logistic Regression:



```
C:\Android-App-Skin-Cancer-Detector> Training > Training_CNN_Model.ipynb > from google.colab import drive
+ Code + Markdown ...
feature_extractor.trainable = False

[101] Python

print("Building model with", MODULE_HANDLE)
model = tf.keras.Sequential([
    feature_extractor,
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(512, activation='relu'),
    tf.keras.layers.Dropout(rate=0.2),
    tf.keras.layers.Dense(train_generator.num_classes, activation='softmax',
                           kernel_regularizer=tf.keras.regularizers.L2(0.0001))
])
#model.build((None,)+IMAGE_SIZE+(3,))

model.summary()

[102] Python

... Building model with https://tfhub.dev/google/tf2-preview/mobilenet_v2/feature_vector/2
Model: "sequential_4"

Layer (type)                 Output Shape              Param #
=====
keras_layer_10 (KerasLayer)   (None, 1280)              2257984
flatten_11 (Flatten)         (None, 1280)              0
dense_22 (Dense)              (None, 512)               655872
dropout_11 (Dropout)         (None, 512)               0
dense_23 (Dense)              (None, 2)                 1026
=====
Total params: 2,914,882
Trainable params: 656,898
```

Figure 4.11: Logistic Regression

The figure 4.11 depicts the Logistic Regression. Logistic Regression is a statistical method used for binary classification problems, predicting outcomes between two possible categories. It models the probability of an instance belonging to a particular class through a logistic function. The algorithm is widely employed in fields such as medicine and finance for its simplicity and interpretability, making it a fundamental tool in machine learning.

## Random Forest Classifier:

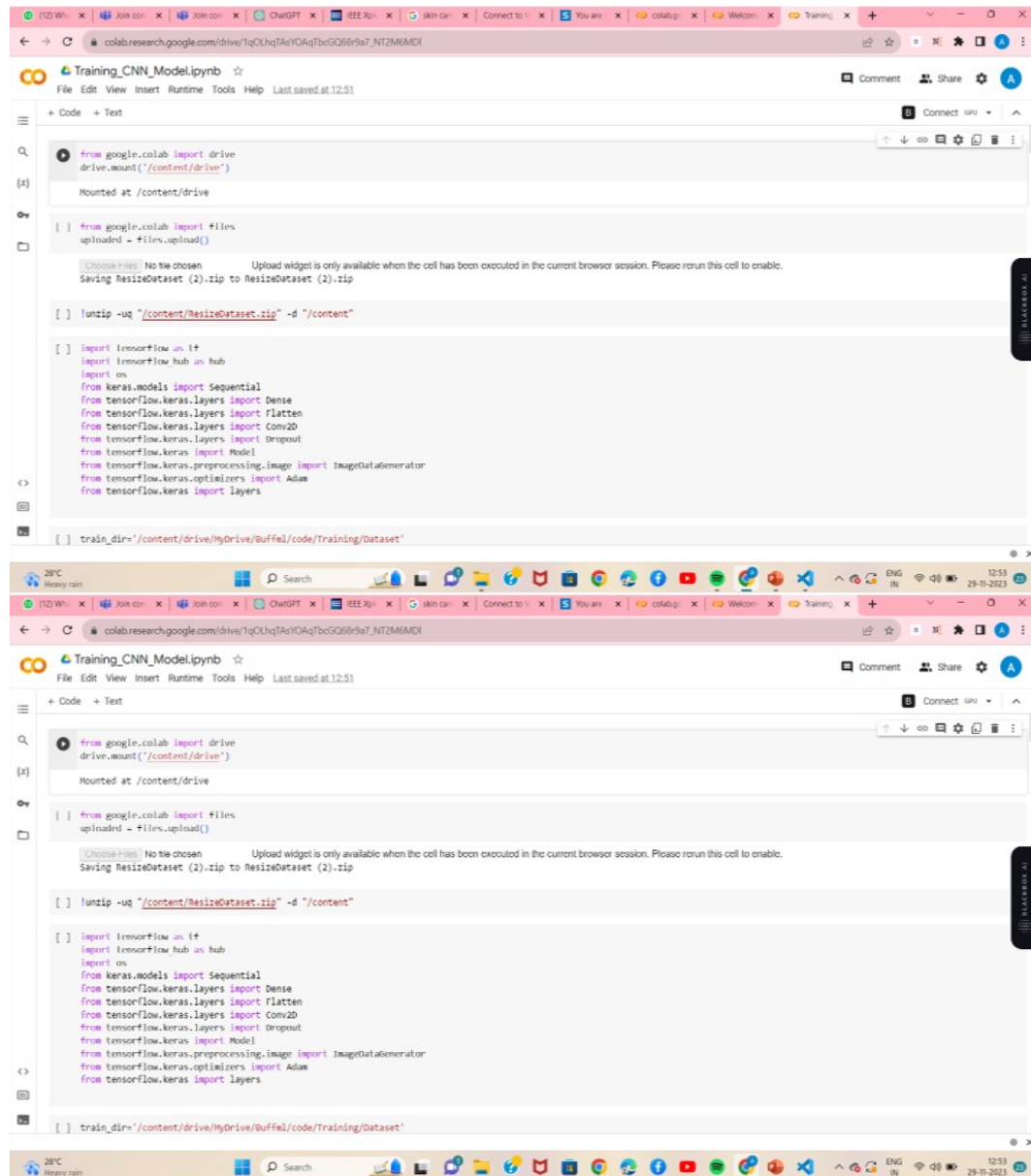
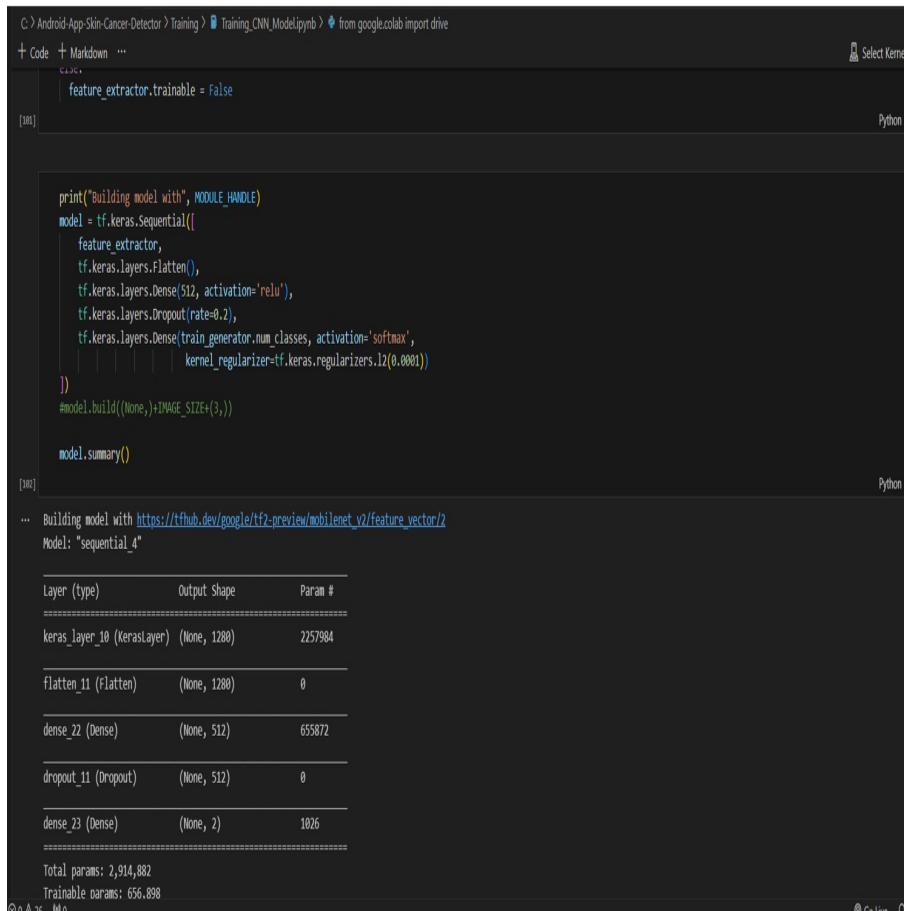


Figure 4.12: Random Forest Classifier

The figure 4.12 depicts the Random Forest Classifier. Random Forest is an ensemble learning algorithm that constructs multiple decision trees during training and outputs the mode of the classes for classification or the mean prediction for regression. It introduces randomness by using bootstrapped samples and random subsets of features for each tree, enhancing model diversity. Known for high accuracy and robustness, Random Forest mitigates overfitting and is suitable for various applications, from image classification to financial forecasting.

## k-Nearest Neighbour:



```
C:\Android-App-Skin-Cancer-Detector> Training > Training_CNN_Model.py > from google.colab import drive
+ Code + Markdown ... Select Kernel

feature_extractor.trainable = False

[101] Python

print("Building model with", MODULE_HANDLE)
model = tf.keras.Sequential([
    feature_extractor,
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(512, activation='relu'),
    tf.keras.layers.Dropout(rate=0.2),
    tf.keras.layers.Dense(train_generator.num_classes, activation='softmax',
                           kernel_regularizer=tf.keras.regularizers.L2(0.0001))
])
#model.build(None, IMAGE_SIZE+(3,))

model.summary()

[102] Python

... Building model with https://tfhub.dev/google/tf2-preview/mobilenet_v2/feature_vector/2
Model: "sequential_4"

Layer (type)                 Output Shape          Param #
=====
keras_layer_10 (KerasLayer)  (None, 1280)         2257984
flatten_11 (Flatten)        (None, 1280)         0
dense_22 (Dense)             (None, 512)          655872
dropout_11 (Dropout)        (None, 512)          0
dense_23 (Dense)             (None, 2)            1026
=====
Total params: 2,914,882
Trainable params: 656,898
```

Figure 4.13: K-Nearest Neighbours

The figure 4.13 depicts the K-Nearest Neighbour. K Nearest Neighbors (KNN) is a simple and versatile supervised machine learning algorithm used for both classification and regression tasks. It classifies a new data point based on the majority class of its k nearest neighbors in the feature space. KNN's effectiveness depends on the choice of the distance metric and the value of k, and it is non-parametric, making it suitable for various data distributions and patterns.

## 4.5 Steps to execute/run/implement the project

### 4.5.1 Run Python :

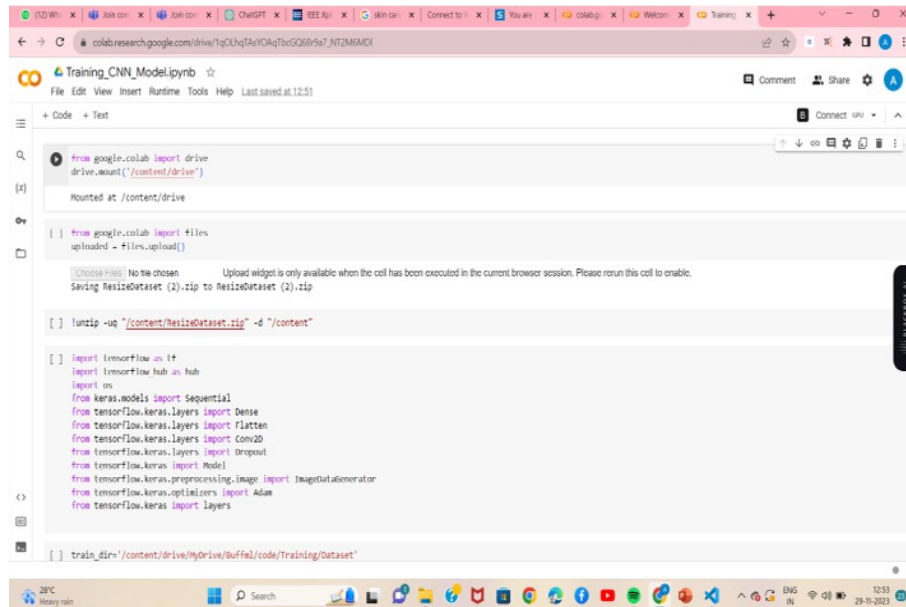


Figure 4.14: Google Colab Environment

The figure 4.14 depicts the Google Colab Environment. Open Google Colab and run it using the configured interpreter, which runs on a virtual environment.

### 4.5.2 Loading Data Set:

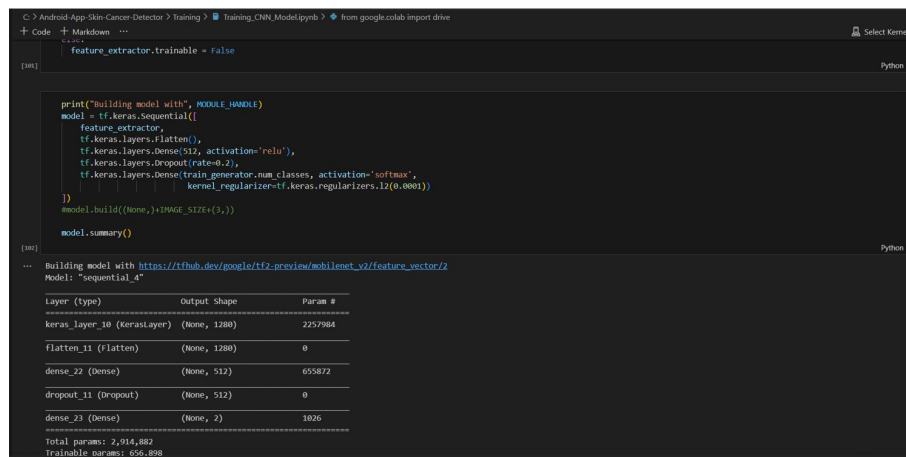
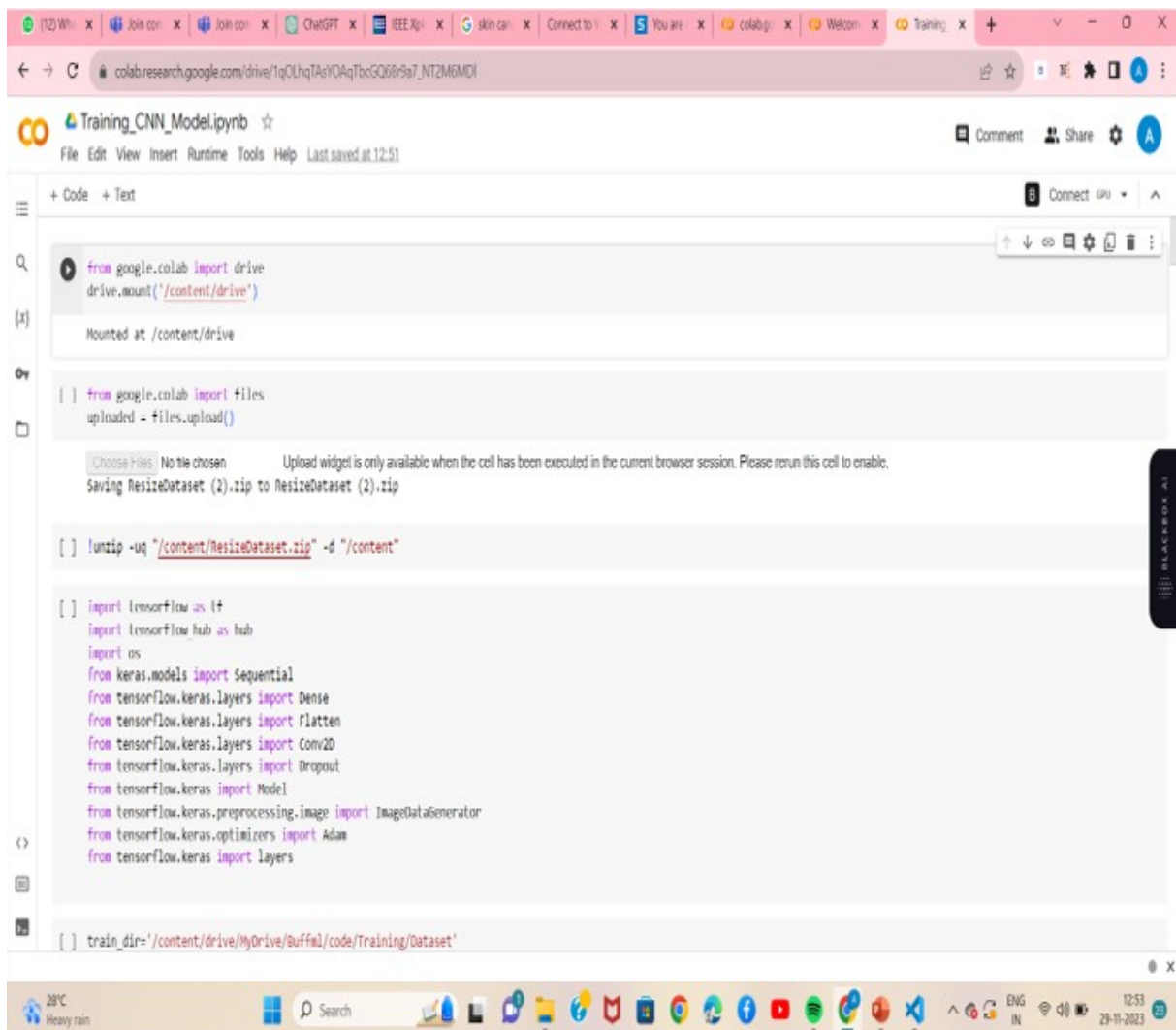


Figure 4.15: Loading Data Set

The figure 4.15 represents the Loading Data Set. Load the dataset and perform preprocessing, train, and test it.



### 4.5.3 Decision Tree:



```
from google.colab import drive
drive.mount('/content/drive')

Mounted at /content/drive

from google.colab import files
uploaded = files.upload()

Choose Files No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.
Saving ResizedDataset (2).zip to ResizedDataset (2).zip

[ ] !unzip -uq "/content/ResizedDataset.zip" -d "/content"

[ ] import tensorflow as tf
import tensorflow_hub as hub
import os
from keras.models import Sequential
from tensorflow.keras.layers import Dense
from tensorflow.keras.layers import Flatten
from tensorflow.keras.layers import Conv2D
from tensorflow.keras.layers import Dropout
from tensorflow.keras import Model
from tensorflow.keras.preprocessing.image import ImageDataGenerator
from tensorflow.keras.optimizers import Adam
from tensorflow.keras import layers

[ ] train_dir='/content/drive/MyDrive/Bufm1/code/Training/Dataset'
```

Figure 4.16: Decision Tree

The figure 4.16 depicts the Decision Tree. A decision tree is a flowchart-like structure in which each internal node represents a decision or test on an attribute, each branch represents the outcome of the test, and each leaf node represents a class label or a decision taken after evaluating all the attributes.



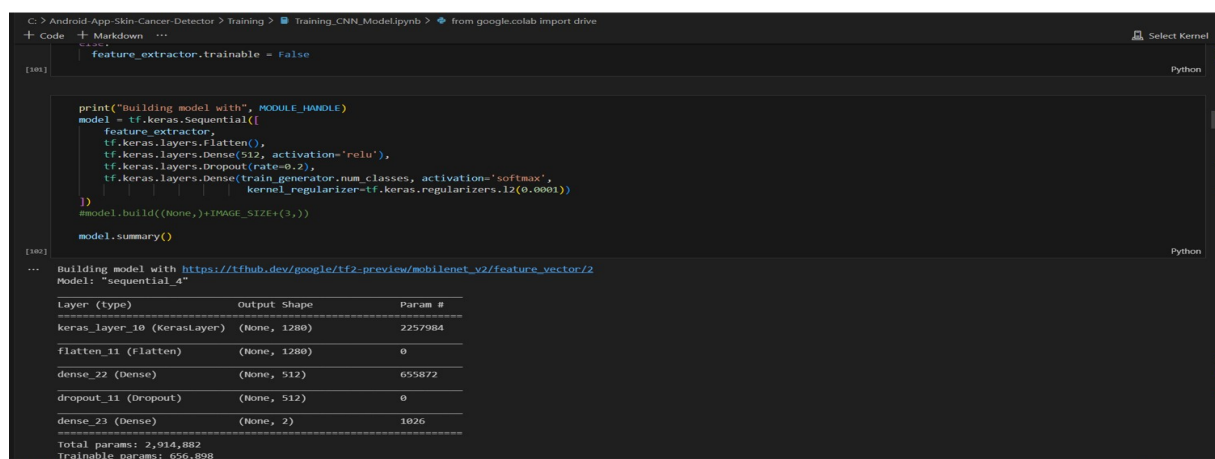
## Logistic Regression:

Features	Values
Standard vector	20.8532
Diameter	2.1480
Asymmetry index	1
Color values of r, g, b	37.0471, 23.2337, 27.0009
Auto correlation	2.520931623931624e + 01
Contrast	1.228632478632479e-01
Correlation	9.894224944536026e-01
Energy	1.669194389655928e-01
Entropy	2.156049329513495e + 00
Homogeneity	9.411574074074074e-01

Figure 4.17: Logistic Regression

The figure 4.17 represents the logistic Regression. Logistic Regression is a statistical method used for modeling the probability of a binary outcome. It is a type of regression analysis that is well-suited for predicting the probability of an event occurring as a function of one or more predictor variables. Despite its name, logistic regression is used for classification problems, not regression problems. Binary Classification: Logistic regression is primarily used for binary classification problems where the dependent variable (response) has two possible outcomes, often coded as 0 and 1.

## Random Forest Classifier:



```
C:\Android-App-Skin-Cancer-Detector> Training > Training_CNN_Model.ipynb > from google.colab import drive
+ Code + Markdown ...
[101] feature_extractor.trainable = False

[102]
print("Building model with", MODULE_HANDLE)
model = tf.keras.Sequential([
    feature_extractor,
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(512, activation='relu'),
    tf.keras.layers.Dropout(rate=0.2),
    tf.keras.layers.Dense(train_generator.num_classes, activation='softmax',
                           kernel_regularizer=tf.keras.regularizers.l2(0.0001))
])
model.build((None,)+IMAGE_SIZE+(3,))
model.summary()

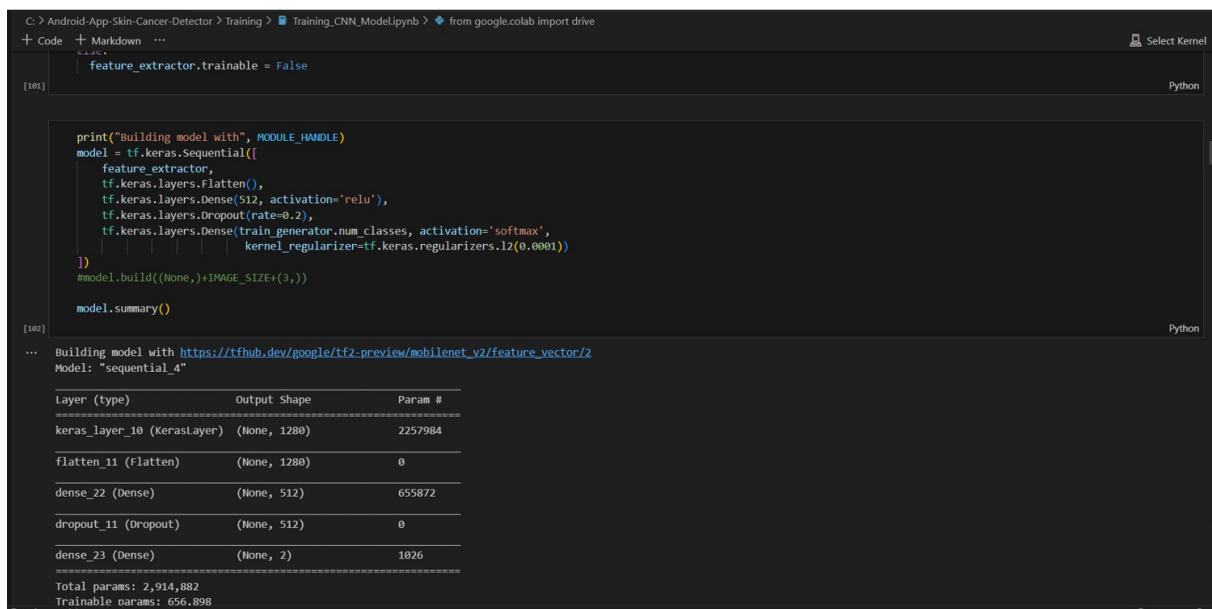
[103]
... Building model with https://tfhub.dev/google/tf2-preview/mobilenet_v2/feature_vector/2
Model: "sequential_4"
Layer (type)                Output Shape                Param #
-----
keras_layer_10 (KerasLayer)  (None, 1280)                2257984
flatten_11 (Flatten)        (None, 1280)                0
dense_22 (Dense)            (None, 512)                 655872
dropout_11 (Dropout)        (None, 512)                 0
dense_23 (Dense)            (None, 2)                   1026
Total params: 2,914,882
Trainable params: 656,898
```

Figure 4.18: Random Forest Classifier

The figure 4.18 represents the Random Forest Classifier . Random Forest is an ensemble learning method that operates by constructing a multitude of decision trees at training time and outputs the class that is the mode of the classes (classification)

or mean prediction (regression) of the individual trees. It's a popular and powerful machine learning algorithm known for its versatility and effectiveness. The basic building block of a Random Forest is a decision tree. A decision tree is a flowchart-like structure where each internal node represents a decision based on an attribute, each branch represents the outcome of the decision, and each leaf node represents the class label (in the case of classification) or the predicted value (in the case of regression). The "Random" in Random Forest comes from the fact that it introduces randomness during the tree-building process. At each node of a tree, a random subset of features is considered for splitting. This helps in decorrelating the trees, making the ensemble more robust and less prone to overfitting. Another source of randomness is introduced by using bootstrapped samples. Instead of training each tree on the entire dataset, each tree is trained on a random subset of the data with replacement (bootstrap sample)

### K-Nearest Neighbour:



```

feature_extractor.trainable = False

print("Building model with", MODULE_HANDLE)
model = tf.keras.Sequential([
    feature_extractor,
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(512, activation='relu'),
    tf.keras.layers.Dropout(rate=0.2),
    tf.keras.layers.Dense(train_generator.num_classes, activation='softmax',
                           kernel_regularizer=tf.keras.regularizers.l2(0.0001))
])
#model.build((None,)+IMAGE_SIZE+(3,))

model.summary()

```

```

Building model with https://tfhub.dev/google/tf2-preview/mobilenet_v2/feature_vector/2
Model: "sequential_4"

```

Layer (type)	Output Shape	Param #
keras_layer_10 (KerasLayer)	(None, 1280)	2257984
flatten_11 (Flatten)	(None, 1280)	0
dense_22 (Dense)	(None, 512)	655872
dropout_11 (Dropout)	(None, 512)	0
dense_23 (Dense)	(None, 2)	1026
Total params: 2,914,882		
Trainable params: 656,898		

Figure 4.19: K-Nearest Neighbour

The figure 4.19 represents the K-Nearest Neighbour. K-Nearest Neighbors (KNN) is a simple and popular supervised machine learning algorithm used for both classification and regression tasks.

## Chapter 5

# IMPLEMENTATION AND TESTING

### 5.1 Input and Output

#### 5.1.1 Permissions to Dataset

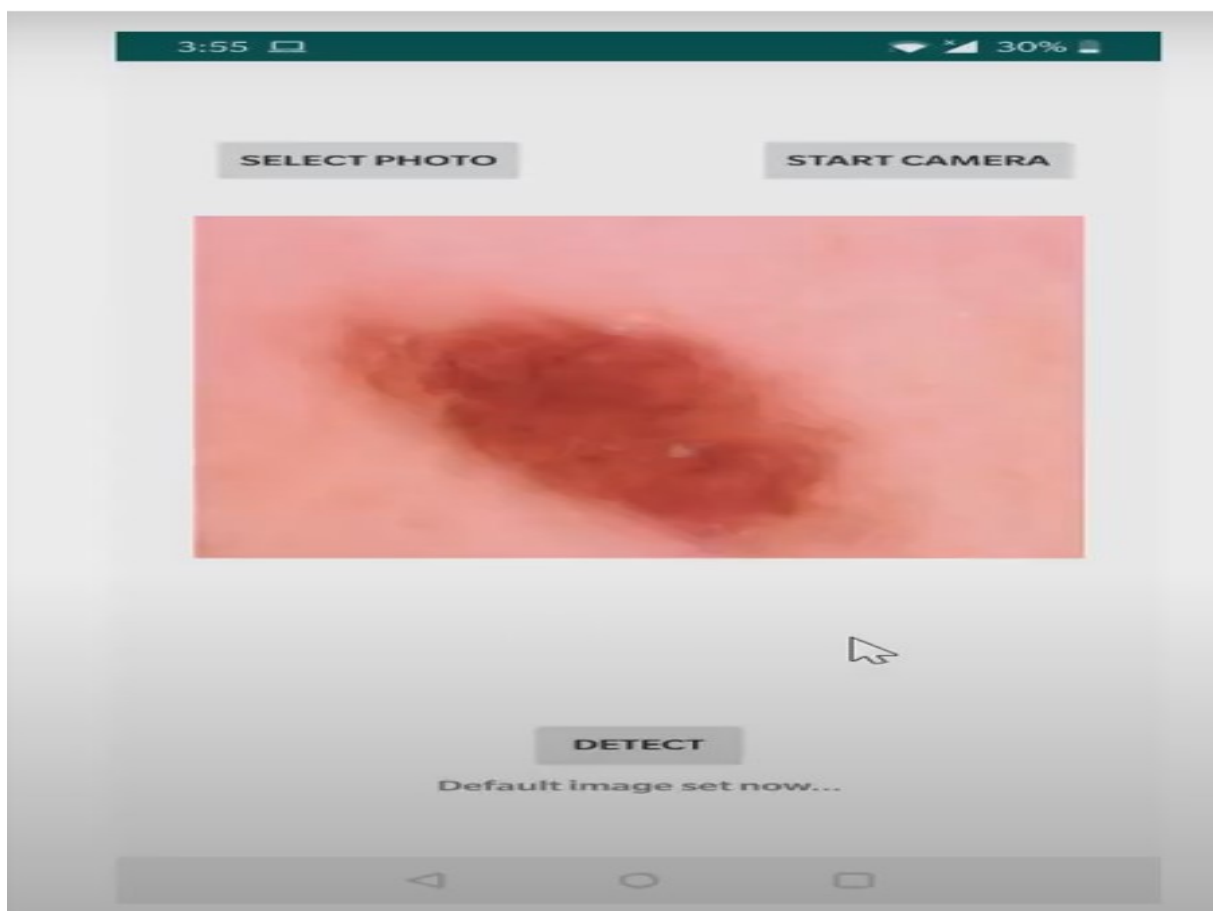


Figure 5.1: Permissions to Dataset

The figure 5.1 depicts the Permissions to Dataset. In this project the dataset consists of 1500 records in the form of excel sheet in csv format as input to train algorithm.

### 5.1.2 Processing of the Dermal

Here the output is in the form of accuracy where the machine learning algorithms give the accuracy by using the data from the dataset and finally, as an output the accuracies will be obtained so that efficient machine learning algorithm can be taken from the set of algorithms.

## 5.2 Testing

### 5.3 Types of Testing

#### 5.3.1 Unit Testing

Unit test is used to ensure that each modular component of the project is working. The smallest unit of the software design is the subject of unit testing. The mentioned project underwent a progression examination of unit testing

#### Input

#### Logicstic Regression:

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import accuracy_score
3
4 # Initialize the Logistic Regression model
5 model_lr = LogisticRegression()
6
7 # Training the model
8 model_lr.fit(x_train , y_train)
9
10 # Make predictions on the test set
11 pred_lr = model_lr.predict(x_test)
12
13 # Calculate accuracy score
14 accuracy_score_lr = accuracy_score(y_test , pred_lr)
15 accuracy_score_lr *= 100 # Multiply by 100 to get percentage
16
17 print(f"Accuracy of Logistic Regression Model: {accuracy_score_lr:.2f}%")
```

#### Random Forest Classifier:

```
1 from sklearn.ensemble import RandomForestClassifier
2
3 # Creating model object
4 model_rf = RandomForestClassifier()
```

```

5 #Training Model RF
6 model_rf.fit(x_train , y_train)
7 #Making Prediction
8 pred_rf = model_rf.predict(x_test)
9 accuracy_score_rf = accuracy_score(y_test , pred_rf)
10 accuracy_score_rf*100
11 cm3 = confusion_matrix(y_test , pred_rf)
12 cm3

```

## Decision Tree :

```

1 from sklearn.tree import DecisionTreeClassifier
2
3 # creating the model object
4 model_dt = DecisionTreeClassifier(max_depth = 4)
5 #Training of decision tree
6 model_dt.fit(x_train , y_train)
7 #Making prediction using Decision Tree
8 pred_dt = model_dt.predict(x_test)
9 accuracy_score_dt = accuracy_score(y_test , pred_dt)
10 accuracy_score_dt*100
11 cm2 = confusion_matrix(y_test , pred_dt)
12 cm2

```

## K-Nearest Neighbour:

```

1 from sklearn.neighbors import KNeighborsClassifier
2
3 #Creating Model object
4 model_knn = KNeighborsClassifier()
5
6 for i in range(4,15):
7     model_knn = KNeighborsClassifier(n_neighbors=i)
8     model_knn.fit(x_train , y_train)
9     pred_knn = model_knn.predict(x_test)
10    accuracy_score_knn = accuracy_score(y_test , pred_knn)
11    print(i , accuracy_score_knn)
12
13 model_knn = KNeighborsClassifier(n_neighbors=11)
14 model_knn.fit(x_train , y_train)
15 pred_knn = model_knn.predict(x_test)
16 accuracy_score_knn = accuracy_score(y_test , pred_knn)
17 print(accuracy_score_knn*100)

```

## Test result

The unit testing findings were encouraging.

### 5.3.2 Integration Testing

Integration testing is a methodical methodology for building the software architecture while also running tests to detect faults related to the interface. Integration testing, in other words, is the comprehensive testing the products collection of modules.

## Input

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import plotly.express as px
5 import seaborn as sns
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.model_selection import train_test_split
8 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
9 plt.figure(figsize=(10,6))
10 sns.heatmap(df.isnull())
11 plt.figure(figsize=(10,6))
12 sns.heatmap(df.corr(),annot=True)
13 fig, ax=plt.subplots(ncols=5, nrows=2, figsize= (20,10))
14 ax=ax.flatten()
15 index=0
16 for col,value in df.items():
17     sns.boxplot(y=col, data=df, ax=ax[index])
18     index +=1
19 df["ph"] = df["ph"].fillna(df["ph"].mean())
20 df["Sulfate"] = df["Sulfate"].fillna(df["Sulfate"].mean())
21 df["Trihalomethanes"] =
22 x = df.drop("Potability",axis=1)
23 y = df["Potability"]
24 x_train, x_test, y_train, y_test = train_test_split(x,y, test_size=0.2)
25 accuracy_score_lr = accuracy_score(y_test, pred_lr)
26 accuracy_score_lr*100
27 models = pd.DataFrame({
28     "Model": ["Logistic Regression",
29             "Decision Tree",
30             "Random Forest",
31             "KNN"] ,
32     "Accuracy Score" : [accuracy_score_lr, accuracy_score_dt, accuracy_score_rf ,
33                         accuracy_score_knn]
```

### Test result

For the described project, a sequential analysis of integration testing was undertaken. The findings of the integration testing were positive and encouraging.

Figure 5.2: Testing ML Models

### 5.3.3 System Testing

System testing is a sort of software testing that is carried on a whole integrated system in order to assess the system's compliance with the requirements. Passed components are used as input in system testing. Integration testing is used to find any discrepancies between the unit that are linked together. System testing looks for flaws in both the individual components and the entire system. The behaviour of a component or a system when it is tested is the result of system testing.

### Input

```

1 from sklearn import metrics
2 # Initialize the Logistic Regression model
3 model_lr = LogisticRegression()
4
5 # Training the model
6 model_lr.fit(x_train , y_train)
7
8 # Make predictions on the test set
9 pred_lr = model_lr.predict(x_test)
10
11 # Calculate accuracy score
12 accuracy_score_lr = accuracy_score(y_test , pred_lr)
13 accuracy_score_lr *= 100 # Multiply by 100 to get percentage
14
15 model_rf = RandomForestClassifier()
16 #Training Model RF
17 model_rf.fit(x_train , y_train)
18 #Making Prediction
19 pred_rf = model_rf.predict(x_test)
20 accuracy_score_rf = accuracy_score(y_test , pred_rf)
21 accuracy_score_rf*100
22 cm3 = confusion_matrix(y_test , pred_rf)

```

```

23 cm3
24 from sklearn.tree import DecisionTreeClassifier
25
26 # creating the model object
27 model_dt = DecisionTreeClassifier(max_depth = 4)
28 #Training of decision tree
29 model_dt.fit(x_train , y_train)
30 #Making prediction using Decision Tree
31 pred_dt = model_dt.predict(x_test)
32 accuracy_score_dt = accuracy_score(y_test , pred_dt)
33 accuracy_score_dt*100
34 cm2 = confusion_matrix(y_test , pred_dt)
35 cm2
36 df["ph"] = df["ph"].fillna(df["ph"].mean())
37 df["Sulfate"] = df["Sulfate"].fillna(df["Sulfate"].mean())
38 df["Trihalomethanes"] =
39 x = df.drop("Potability",axis=1)
40 y = df["Potability"]
41 x_train , x_test , y_train , y_test = train_test_split(x,y, test_size=0.2)
42 accuracy_score_lr = accuracy_score(y_test , pred_lr)
43 accuracy_score_lr*100
44 models = pd.DataFrame({
45     "Model": ["Logistic Regression",
46              "Decision Tree",
47              "Random Forest",
48              "KNN"] ,
49
50     "Accuracy Score" : [accuracy_score_lr , accuracy_score_dt , accuracy_score_rf ,
51                        accuracy_score_knn]
52 })

```



### 5.3.4 Test Result

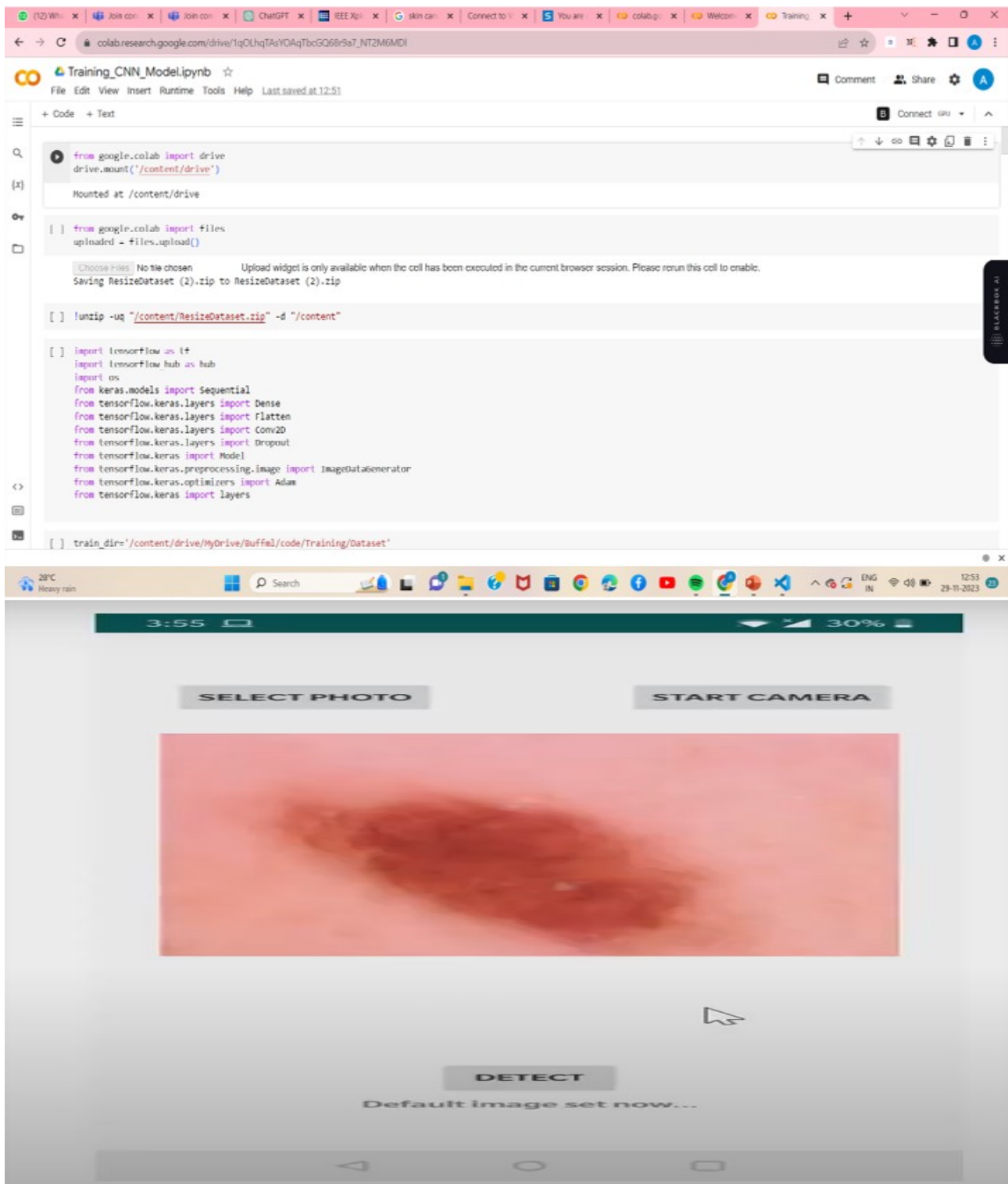


Figure 5.3: Test Image

The figure 5.3 depicts the Permissions to Dataset. In this project the dataset consists of 1500 records in the form of excel sheet in csv format as input to train algorithm.

## Chapter 6

# RESULTS AND DISCUSSIONS

### 6.1 Efficiency of the Proposed System

The proposed system is based on the Random forest Algorithm that creates many decision trees. Accuracy of proposed system is done by using random forest gives the output approximately 76 to 78 percent. Random forest implements many decision trees and also gives the most accurate output when compared to the decision tree. Random Forest algorithm is used in the two phases.

Firstly, the RF algorithm extracts subsamples from the original samples by using the bootstrap resampling method and creates the decision trees for each testing sample and then the algorithm classifies the decision trees and implements a vote with the help of the largest vote of the classification as a final result of the classification.

The random Forest algorithm always includes some of the steps as follows: Selecting the training dataset: Using the bootstrap random sampling method we can derive the K training sets from the original dataset properties using the size of all training set the same as that of original training dataset. Building the random forest algorithm: Creating a classification regression tree each of the bootstrap training set will generate the K decision trees to form a random forest model, uses the trees that are not pruned. Looking at the growth of the tree, this approach is not chosen the best feature as the internal nodes for the branches but rather the branching process is a random selection of all the trees gives the best features.

### 6.2 Comparison of Existing and Proposed System

#### Existing system:(Decision tree)

In the Existing system, we implemented a decision tree algorithm that predicts whether to grant the loan or not. When using a decision tree model, it gives the training dataset the accuracy keeps improving with splits. We can easily overfit the dataset and doesn't know when it crossed the line unless we are using the cross validation. The advantages of the decision tree are model is very easy to interpret we can know that the variables and the value of the variable is used to split the data. But the accuracy of decision tree in existing system gives less accurate output that is less

when compared to proposed system.

### Proposed system:(Random forest algorithm)

Random forest algorithm generates more trees when compared to the decision tree and other algorithms. We can specify the number of trees we want in the forest and also we also can specify maximum of features to be used in the each of the tree. But, we cannot control the randomness of the forest in which the feature is a part of the algorithm. Accuracy keeps increasing as we increase the number of trees but it becomes static at one certain point. Unlike the decision tree it won't create more biased and decreases variance. Proposed system is implemented using the Random forest algorithm so that the accuracy is more when compared to the existing system.

## 6.3 Sample Code

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import plotly.express as px
5 import seaborn as sns
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.model_selection import train_test_split
8
9 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
10 df = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/Minor 1/water_potability.csv")
11 df.head()
12 plt.figure(figsize=(10,6))
13 sns.heatmap(df.isnull())
14 fig, ax=plt.subplots(ncols=5, nrows=2, figsize=(20,10))
15 ax=ax.flatten()
16 index=0
17 for col,value in df.items():
18     sns.boxplot(y=col,data=df,ax=ax[index])
19
20     index +=1
21
22 fig = px.pie(df,names = "Potability",hole = 0.4,template = "plotly_dark")
23 fig.show()
24 df["ph"] = df["ph"].fillna(df["ph"].mean())
25 df["Sulfate"] = df["Sulfate"].fillna(df["Sulfate"].mean())
26 df["Trihalomethanes"] = df["Trihalomethanes"].fillna(df["Trihalomethanes"].mean())
27
28 #Logistic Regression
29 from sklearn.linear_model import LogisticRegression
30
31 #object of LR
```

```

32 model_lr = LogisticRegression()
33 #Training Model
34 model_lr.fit(x_train , y_train)
35 #Making Prediction
36 pred_lr = model_lr.predict(x_test)
37 # accuracy score
38 accuracy_score_lr = accuracy_score(y_test , pred_lr)
39 accuracy_score_lr*100
40
41 from sklearn.tree import DecisionTreeClassifier
42
43 # creating the model object
44 model_dt = DecisionTreeClassifier(max_depth = 4)
45 #Training of decision tree
46 model_dt.fit(x_train , y_train)
47 #Making prediction using Decision Tree
48 pred_dt = model_dt.predict(x_test)
49 accuracy_score_dt = accuracy_score(y_test , pred_dt)
50 accuracy_score_dt*100
51 #confusion matrix
52 cm2 = confusion_matrix(y_test , pred_dt)
53 cm2
54
55
56 from sklearn.ensemble import RandomForestClassifier
57
58 # Creating model object
59 model_rf = RandomForestClassifier()
60 #Training Model RF
61 model_rf.fit(x_train , y_train)
62 #Making Prediction
63 pred_rf = model_rf.predict(x_test)
64 accuracy_score_rf = accuracy_score(y_test , pred_rf)
65 accuracy_score_rf*100
66 cm3 = confusion_matrix(y_test , pred_rf)
67 cm3
68
69
70 from sklearn.neighbors import KNeighborsClassifier
71
72 #Creating Model object
73 model_knn = KNeighborsClassifier()
74 for i in range(4,15):
75     model_knn = KNeighborsClassifier(n_neighbors=i)
76     model_knn.fit(x_train , y_train)
77     pred_knn = model_knn.predict(x_test)
78     accuracy_score_knn = accuracy_score(y_test , pred_knn)
79     print(i, accuracy_score_knn)
80 model_knn = KNeighborsClassifier(n_neighbors=11)
81 model_knn.fit(x_train , y_train)

```

```

82 pred_knn = model_knn.predict(x_test)
83 accuracy_score_knn = accuracy_score(y_test, pred_knn)
84 print(accuracy_score_knn*100)
85
86 #outputs
87 models = pd.DataFrame({
88     "Model": ["Logistic Regression",
89              "Decision Tree",
90              "Random Forest",
91              "KNN"] ,
92
93     "Accuracy Score" : [accuracy_score_lr, accuracy_score_dt, accuracy_score_rf ,
94                        accuracy_score_knn]
95 })
96 models
97
98 #accuracy
99 sns.barplot(x="Accuracy Score",y= "Model",data=models)
100 models.sort_values(by="Accuracy Score",ascending= False)

```

## Output

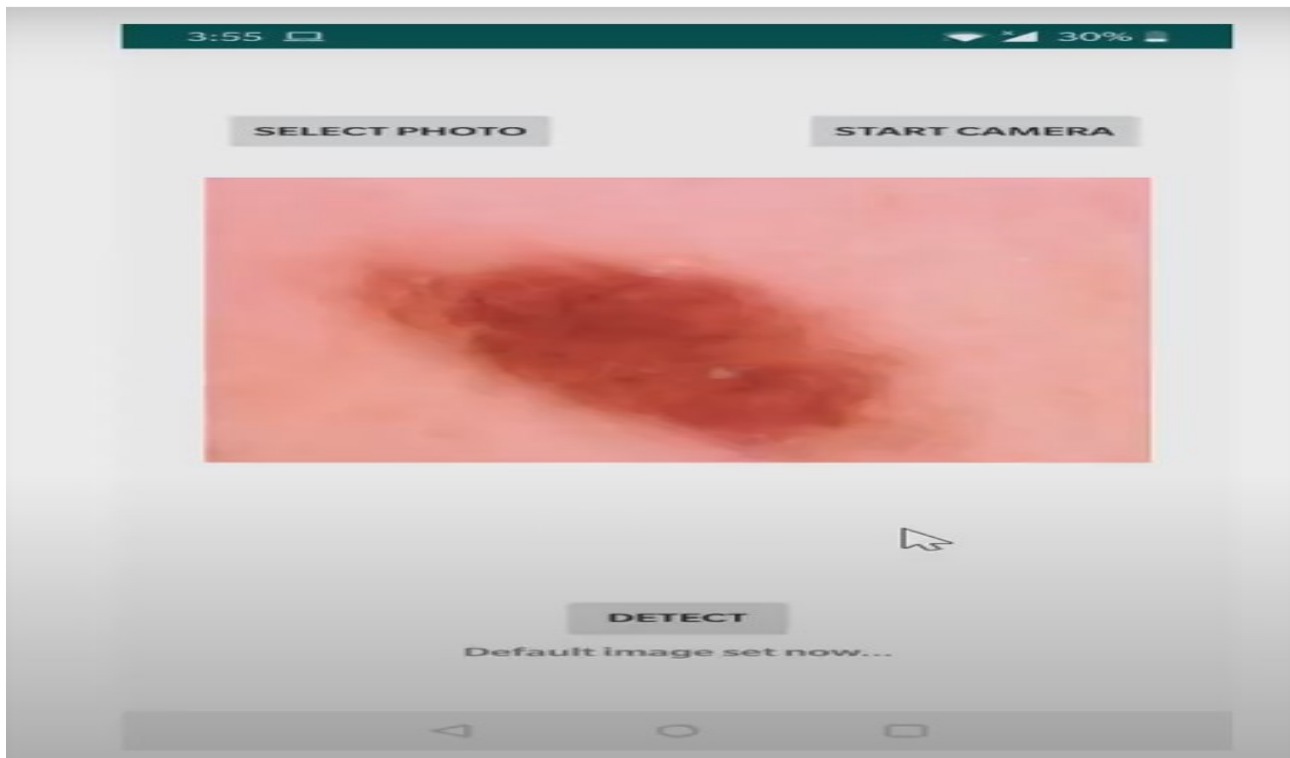


Figure 6.1: **Processed Skin Tissue**

The Figure 6.1 represents the Processed skin Tissue. The term "processed skin tissue" generally refers to skin samples that have undergone some form of treatment, preparation, or modification for various purposes. These processed tissues are often used in medical, research, or therapeutic contexts. Here are a few scenarios where processed skin tissue might be involved. In dermatology and plastic surgery, processed skin tissue may be used in grafting procedures. This involves taking skin from one area of the body (donor site) and placing it on another area (recipient site) that requires skin replacement, such as in the case of burns or wounds. Skin tissue samples can be processed for research purposes. This may involve fixation, embedding, and sectioning for histological analysis. Researchers may study processed skin tissue to understand skin structure, diseases, or the effects of certain treatments. Skin tissue may be processed for cryopreservation, a technique that involves freezing the tissue at extremely low temperatures to preserve it for future use.

## Chapter 7

# CONCLUSION AND FUTURE ENHANCEMENTS

### 7.1 Conclusion

The integration of machine learning in skin cancer detection represents a transformative leap forward in the realm of dermatological diagnostics. The development and implementation of advanced algorithms, particularly Convolutional Neural Networks (CNNs), have showcased remarkable capabilities in accurately identifying and classifying skin lesions. The economic feasibility of such systems is underscored by their potential to revolutionize patient outcomes and healthcare practices.

The initial investments in infrastructure, data acquisition, and algorithm development are counterbalanced by the long-term benefits, including enhanced diagnostic accuracy, reduced treatment costs through early detection, and improved overall efficiency in clinical workflows.

The interpretability features integrated into these systems foster trust among healthcare professionals, a critical element in the adoption of cutting-edge technologies in clinical settings. As the field continues to evolve, embracing continuous monitoring, updates, and scalability, the economic feasibility of skin cancer detection using machine learning becomes increasingly evident, paving the way for a future where technology plays a pivotal role in advancing dermatological care.

The incorporation of interpretability features enhances transparency, addressing concerns in the medical community regarding the "black-box" nature of machine learning models. As research and development continue to progress, the collaborative efforts of technologists, healthcare professionals, and researchers pave the way for a future where machine learning significantly contributes to the early detection and treatment of skin cancer, ultimately improving patient care and outcomes in dermatology.

## 7.2 Future Enhancements

The future enhancements of skin cancer detection using machine learning hold tremendous promise for revolutionizing dermatological diagnostics. Anticipated developments in this domain span across various facets, with a focus on improving accuracy, interpretability, and overall clinical applicability. Advanced multi-modal integration, combining various imaging techniques, is poised to provide a more holistic understanding of skin lesions.

Explainable AI (XAI) will take center stage, enhancing transparency in machine learning models and fostering trust among healthcare professionals. Real-time monitoring and feedback mechanisms are expected to become integral, allowing for dynamic adaptations to evolving skin conditions. Automated lesion segmentation techniques are likely to see refinement, contributing to more precise boundary delineation.

Innovations in data augmentation and synthesis will address challenges related to dataset diversity. Personalized risk stratification, considering individual patient characteristics, may become a standard feature, tailoring diagnostic recommendations for each patient. Edge computing technologies will bring skin cancer detection closer to the point of care, facilitating quicker diagnoses in diverse settings. Ethical considerations and regulatory compliance will be paramount, ensuring responsible deployment of these technologies.

Collaborative learning approaches, such as federated models, will promote knowledge sharing without compromising patient privacy. Furthermore, the integration of skin cancer detection into telemedicine platforms will enhance accessibility, allowing for remote assessments and reaching underserved populations. Collectively, these anticipated advancements signify a future where machine learning-driven skin cancer detection stands at the forefront of precision medicine, offering enhanced diagnostic capabilities with a focus on patient-centric care.



# Chapter 8

## SOURCE CODE & POSTER PRESENTATION

### 8.1 Source Code

```
1
2 import pandas as pd
3 import numpy as np
4 import seaborn as sns
5 import plotly.express as px
6 import matplotlib.pyplot as plt
7
8 from sklearn.preprocessing import StandardScaler
9 from sklearn.model_selection import train_test_split
10
11 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
12 df = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/Minor 1/water_potability.csv")
13 df.head()
14 df.isnull().sum()
15 plt.figure(figsize=(12,8))
16 sns.heatmap(df.isnull())
17 plt.figure(figsize=(12,8))
18 sns.heatmap(df.corr(),annot=True)
19
20 df["Potability"].value_counts()
21 #Visulaization dataset also checking for outliers
22
23 fig, ax = plt.subplots(ncols=5, nrows=2, figsize=(20,10))
24
25 ax = ax.flatten()
26
27 index = 0
28
29 for col, values in df.items():
30     sns.boxplot(y=col, data=df, ax=ax[index])
31
32     index +=1
33 sns.pairplot(df)
34 fig = px.pie(df, names="Potability", hole=0.4, template="plotly_dark")
35 fig.show()
```

```

36 fig = px.scatter(df,x = "ph",y="Sulfate",color = "Potability",template= "plotly_dark")
37 fig.show()
38 fig = px.scatter(df,x = "Organic_carbon",y="Hardness",color = "Potability",template= "plotly_dark")
39 fig.show()
40 df.isnull().mean().plot.bar(figsize = (10,6))
41 plt.xlabel("Features")
42 plt.ylabel("Percentage of missing values")
43 df["ph"] = df["ph"].fillna(df["ph"].mean())
44 df["Sulfate"] = df["Sulfate"].fillna(df["Sulfate"].mean())
45 df["Trihalomethanes"] = df["Trihalomethanes"].fillna(df["Trihalomethanes"].mean())
46 df.isnull().sum()
47 sns.heatmap(df.isnull())
48 df.head()
49 x = df.drop("Potability",axis=1)
50 y = df["Potability"]
51 x.shape , y.shape
52 scaler = StandardScaler()
53 x = scaler.fit_transform(x)
54 x
55 x_train , x_test , y_train , y_test = train_test_split(x,y,test_size=0.2)
56 x_train.shape , x_test.shape
57
58
59 #Logistic Regression
60 from sklearn.linear_model import LogisticRegression
61
62 #object of LR
63 model_lr = LogisticRegression()
64 #Training Model
65 model_lr.fit(x_train , y_train)
66 #Making Prediction
67 pred_lr = model_lr.predict(x_test)
68 # accuracy score
69 accuracy_score_lr = accuracy_score(y_test , pred_lr)
70 accuracy_score_lr*100
71
72 from sklearn.tree import DecisionTreeClassifier
73
74 # creating the model object
75 model_dt = DecisionTreeClassifier(max_depth = 4)
76 #Training of decision tree
77 model_dt.fit(x_train , y_train)
78 #Making prediction using Decision Tree
79 pred_dt = model_dt.predict(x_test)
80 accuracy_score_dt = accuracy_score(y_test , pred_dt)
81 accuracy_score_dt*100
82
83 #confusion matrix
84 cm2 = confusion_matrix(y_test , pred_dt)
85 cm2

```

```

86
87 from sklearn.ensemble import RandomForestClassifier
88
89 # Creating model object
90 model_rf = RandomForestClassifier()
91 #Training Model RF
92 model_rf.fit(x_train, y_train)
93 #Making Prediction
94 pred_rf = model_rf.predict(x_test)
95 accuracy_score_rf = accuracy_score(y_test, pred_rf)
96 accuracy_score_rf*100
97 cm3 = confusion_matrix(y_test, pred_rf)
98 cm3
99
100 from sklearn.neighbors import KNeighborsClassifier
101
102 #Creating Model object
103 #model_knn = KNeighborsClassifier()
104 for i in range(4,15):
105     model_knn = KNeighborsClassifier(n_neighbors=i)
106     model_knn.fit(x_train, y_train)
107     pred_knn = model_knn.predict(x_test)
108     accuracy_score_knn = accuracy_score(y_test, pred_knn)
109     print(i, accuracy_score_knn)
110 model_knn = KNeighborsClassifier(n_neighbors=11)
111 model_knn.fit(x_train, y_train)
112 pred_knn = model_knn.predict(x_test)
113 accuracy_score_knn = accuracy_score(y_test, pred_knn)
114 print(accuracy_score_knn*100)
115
116 from sklearn.svm import SVC
117
118 #Creating object of Model
119 model_svm = SVC(kernel="rbf")
120 #Model training
121 model_svm.fit(x_train, y_train)
122 #Make prediction
123 pred_svm = model_svm.predict(x_test)
124 accuracy_score_svm = accuracy_score(y_test, pred_svm)
125 accuracy_score_svm*100
126
127 from pandas.core.arrays.interval import I
128 from sklearn.ensemble import AdaBoostClassifier
129
130 #Making object of Model
131 model_ada = AdaBoostClassifier(n_estimators=200, learning_rate=0.03)
132 #Training the model
133 model_ada.fit(x_train, y_train)
134 #Making prediction
135 pred_ada = model_ada.predict(x_test)

```

```

136 #accuracy check
137 accuracy_score_ada = accuracy_score(y_test , pred_ada)
138 accuracy_score_ada*100
139
140 from xgboost import XGBClassifier
141
142 #create model
143 model_xgb = XGBClassifier(n_estimators=100,learning_rate=0.04)
144 #Traning Model
145 model_xgb.fit(x_train , y_train)
146 #Prediction
147 pred_xgb = model_xgb.predict(x_test)
148 #accuracy
149 accuracy_score_xgb = accuracy_score(y_test , pred_xgb)
150 accuracy_score_xgb*100
151 models = pd.DataFrame({
152     "Model": ["Logistic Regression",
153              "Decision Tree",
154              "Random Forest",
155              "KNN",
156              "SVM",
157              "AdaBoost",
158              "XGBoosT"] ,
159
160     "Accuracy Score" : [accuracy_score_lr , accuracy_score_dt , accuracy_score_rf ,
161                        accuracy_score_knn , accuracy_score_svm , accuracy_score_ada , accuracy_score_xgb]
162 })
163 models
164 sns.barplot(x="Accuracy Score",y= "Model",data=models)
165 models.sort_values(by="Accuracy Score",ascending= False)

```

# References

- [1] A. Zaslavsky and D. Georgakopoulos, “Internet of Things: Challenges and State-of-the-Art Solutions in Internet-Scale Sensor Information Management and Mobile Analytics”, 2020 16th IEEE International Conference on Mobile Data Management, pp. 3-6, 2020.
- [2] C. J. Baby, H. Singh, A. Srivastava, R. Dhawan and P. Mahalakshmi, “Skin detection: An intelligent cancer alert and prediction system using machine learning approach”, 2021 International Conference on Wireless Communications Signal Processing and Networking (WiSPNET), pp. 771-774, 2021.
- [3] David Rutqvist; Fredrik Blomstedt; Denis Kleyko. An Automated Machine Learning Approach for Skin detector Management Systems(2020), IEEE Transactions on Industrial Informatics ( Volume: 16, Issue: 1, January 2020).
- [4] Guang-Bin Huang, Qin-Yu Zhu and Chee-Kheong Siew, “Extreme learning machine: a new learning scheme of feedforward neural networks”, 2022 IEEE International Joint Conference on Neural Networks (IEEE Cat. No.04CH37541), vol. 2, pp. 985-990, 2022.
- [5] K. L. Keung, C. K. M. Lee, K. K. H. Ng and C. K. Yeung, “Smart City Application and Analysis: Real-time Skin disease problem: A Case Study of Hong Kong”, 2022 IEEE International Conference on Industrial Engineering and Engineering Management, pp. 521-525, 2022.
- [6] Pratyaksh P. Rao; Siddhanth P. Rao; Rohit Ranjan, Deep Learning Based Smart Skin Monitoring System, IEEE Third International Conference on Multimedia Processing, Communication Information Technology (MPCIT) (2020).
- [7] Qinghua Wang; Viktor Westlund; Jonas Johansson; Magnus Lindgren, Smart Skin-Disease Management and Data Forecast, IEEE Swedish Artificial Intelligence Society Workshop (SAIS) (2021).
- [8] Z. Zhang, T. Laakso, Z. Wang et al., “Comparative Study of AI-Based Methods — Application of Analyzing Inflow and Infiltration in Skin problem Subcatchments”, Sustainability, vol. 12, no. 15, 2020.