

AI-Powered Medical Prediction System

A Project Report

submitted in partial fulfillment of the requirements

of

AICTE Internship on AI: Transformative Learning
with

TechSaksham – A joint CSR initiative of Microsoft & SAP

by

Aniket

aniket.kr2103@gmail.com

Under the Guidance of

Pavan Sumohana

ACKNOWLEDGEMENT

I would like to convey my sincere thanks to my mentor, Pavan Sumohana, for their continuous support, advice, and motivation during this project. Their experience, knowledge, and valuable suggestions played a significant role in shaping my knowledge about AI-based medical prediction systems. Their guidance not only assisted me in overcoming difficulties but also helped me execute the project successfully with enhanced accuracy and efficiency.

I am also deeply thankful to AICTE TechSaksham, Microsoft, and SAP for facilitating a platform for collaborative learning. This internship has had a great impact on my technical skills by exposing me to practical applications of AI in healthcare. The experience has expanded my knowledge and bolstered my capability for developing and deploying AI models.

Lastly, I wish to express my heartfelt gratitude towards my friends, family, and colleagues for the encouragement and motivation they provided. Their encouragement sustained me throughout the project, helping me to fulfill this project with success.

ABSTRACT

The rising incidence of chronic diseases is a major challenge to healthcare systems worldwide. Proper and timely diagnosis can improve the effectiveness of treatment and patient outcomes. This project utilizes machine learning methods to forecast disease probabilities from patient information, allowing for early diagnosis and preventive measures.

The process consists of five major steps:

1. Data Collection & Integration – Collating medical datasets from trusted sources, such as publicly available health data repositories (e.g., Kaggle).
2. Data Preprocessing & Cleaning – Maintaining consistency, missing value handling, and data format standardization for enhanced model accuracy.
3. Exploratory Data Analysis (EDA) – Discovering trends, feature correlations, and data distribution patterns to optimize model performance.
4. Machine Learning Model Building – Using Random Forest Classifier and various algorithms for the prediction of diseases. Separate models are trained for Diabetes, Heart Disease, Kidney Disease, Liver Disease, Breast Cancer, Lung Cancer, and Parkinson's Disease.
5. Deployment & Visualization – Developing an interactive web-based interface with Streamlit and Python libraries for real-time prediction of diseases.

This research effectively proves the efficacy of AI-based disease prediction. Future enhancements involve incorporating real-time data processing and extending the model to include more diseases, making it more diagnostic in nature.



TABLE OF CONTENT

Abstract.....
Chapter1.Introduction.....
1.1 Problem Statement
1.2 Motivation
1.3 Objectives
1.4 Scope of the Project
Chapter 2. Literature Survey.....
2.1 Review of Relevant Literature
2.2 Machine Learning Concepts: Random Forest Classifier
2.3 Gaps in Existing Solutions
Chapter 3. Proposed Methodology.....
3.1 System Design
3.2 Requirements Specification
3.2.1 Hardware Requirements
3.2.2 Software Requirements
Chapter 4. Implementation and Results.....
4.1 Model Development for Disease Prediction
4.2 Web App Development using Streamlit
4.3 Snapshots of Results
4.4 GitHub Link for the Code
Chapter 5. Discussion .. and Conclusion.....
5.1 Future Work
5.2 Conclusion
References.....

LIST OF FIGURES

Figure No.	Figure Caption	Page No.
Figure 1	Data Preparation Process	3
Figure 2	Diabetes Predictive Model	5
Figure 3	Heart Disease Predictive Model	5
Figure 4	Kidney Disease Predictive Model	6
Figure 5	Liver Disease Prediction Results	6
Figure 6	Breast Cancer Prediction Results	7
Figure 7	Lung Cancer Prediction Results	8
Figure 8	Parkinson's Disease Prediction Results	8
Figure 9	Streamlit Web Application UI	8



CHAPTER 1

Introduction

1.1 Problem Statement:

Healthcare facilities gather enormous quantities of patient data on a daily basis, but they struggle to analyze it efficiently for predicting diseases at an early stage. Existing conventional diagnostic tools are not very accurate and are susceptible to human errors, and therefore, cause delayed or wrong diagnoses. The lack of predictive systems based on AI technology restricts the efficacy of early interventions, promoting health risk and treatment expenses.

1.2 Motivation:

The objective of this project is to harness the potential of machine learning for creating an automated disease prediction system. By incorporating AI models with real-time patient information, the project hopes to aid healthcare professionals in making timely and accurate diagnoses. This will help initiate early interventions, lower rates of mortality, and improve the efficiency of patient care.

1.3 Objectives:

- Apply machine learning models (Random Forest Classifier, Decision Tree, etc.) for predicting diseases.
- Create an interactive web-based dashboard based on Streamlit for live predictions.
- Automate prediction and data processing workflows to facilitate streamlined healthcare decision-making.
- Enhance diagnostic accuracy and enable clinical decision-making with insights powered by AI.

1.4 Project Scope:

The project targets the utilization of Python and machine learning methods to make predictions for several diseases, namely Diabetes, Heart Disease, Kidney Disease, Liver Disease, Breast Cancer, Lung Cancer, and Parkinson's Disease.

The solution entails:

- Data processing: Patient data cleaning, standardization, and preparation.
- Model training: Utilizing multiple machine learning models, all trained on disease-specific data.
- Visualization: Presenting the prediction results in an easily understandable interface using Streamlit.

Although the project is insightful, prediction accuracy is subjective to data quality, model accuracy, and testing in the real world.

CHAPTER 2

Literature Survey

2.1 Review of Relevant Literature

Current research highlights the increasing importance of machine learning algorithms in healthcare for predictive analytics and early disease detection. Models such as Random Forest, Decision Trees, Support Vector Machines (SVM), and Logistic Regression have been used for disease classification and risk prediction.

Methods like ensemble learning and deep learning models (neural networks) have demonstrated enhanced accuracy in medical diagnosis. Real-time prediction and easy deployment, though, are still topics of research.

2.2 Current Models, Methods

The project deploys Random Forest Classifier and other machine learning algorithms based on Python libraries like NumPy, Pandas, Scikit-learn, and Joblib.

The models are trained from disease-specific data and employ feature engineering and hyperparameter optimization to improve accuracy.

The project examines patient health information to forecast the probability of diseases like Diabetes, Heart Disease, Kidney Disease, Liver Disease, Breast Cancer, Lung Cancer, and Parkinson's Disease, making real-time forecasts.

2.3 Gaps in Current Solutions

There are still gaps in current solutions in areas like:

- Model Accuracy and Reliability: Current models tend to lack accuracy because of unbalanced or missing datasets.
- Real-Time Integration: Most solutions are not based on real-time prediction, rendering them less clinically applicable.
- Data Privacy and Security: Safeguarding sensitive patient information is a persistent concern with healthcare AI deployments.

CHAPTER 3

Proposed Methodology

3.1 System Design

The system has the following primary phases:

- Data Collection: Collecting disease-specific data sets from authentic sources like Kaggle and UCI Machine Learning Repository.
- Data Preprocessing & Feature Engineering: Preprocessing the data set by managing missing values, converting categorical data into numerical data, and feature scaling for numerical features.
- Model Development: Developing individual machine learning models (Random Forest, Decision Tree, Logistic Regression, etc.) for every disease.
- Visualization & Deployment:
 - Developing an interactive Streamlit web app for real-time disease prediction.
 - Presenting the results of prediction and performance indicators in a way that is easily understandable by users.



3.2 Requirement Specification

Listing down the tools and technologies needed to deploy the solution.

3.2.1 Hardware Requirements:

- Processor: AMD Ryzen 5 / Intel i5 or greater
- RAM: 8GB or greater
- Storage: 1GB or greater
- Operating System: Windows 10 / Linux / macOS

3.2.2 Software Requirements:

- Programming Language: Python

Libraries:

- Data Processing: Pandas, NumPy
- Machine Learning: Scikit-learn, Joblib
- Web Development: Streamlit
- Development Environment: VS Code, Jupyter Notebook

CHAPTER 4

Implementation and Result

4.1 Model Building for Disease Prediction

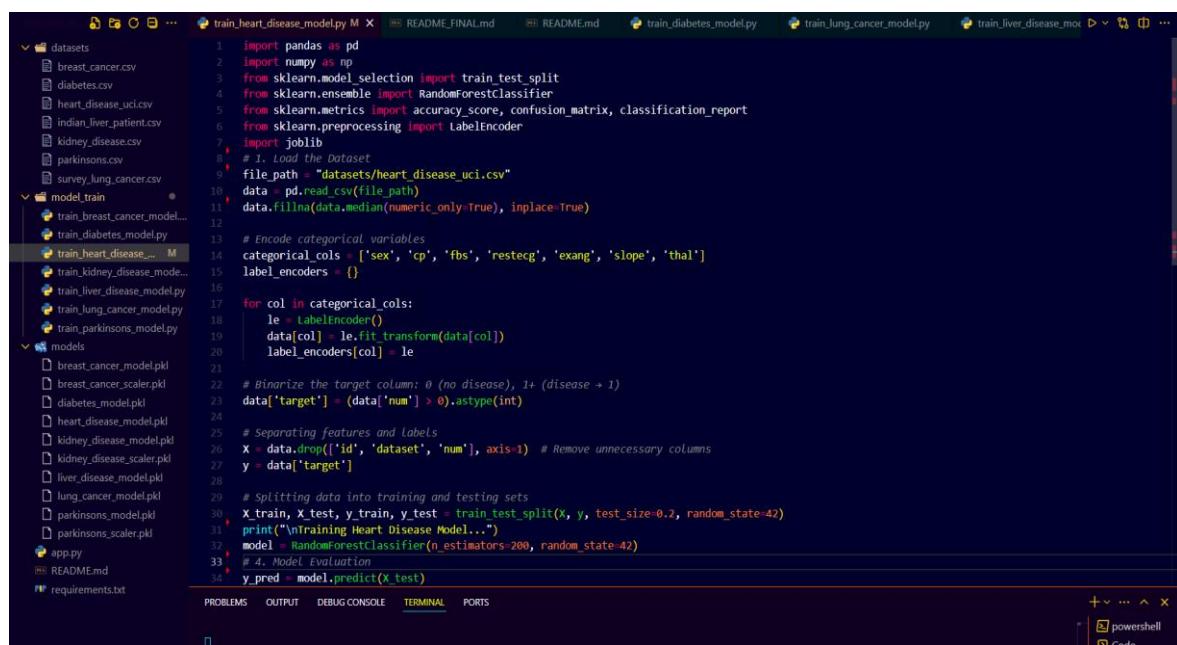
Developing several prediction models for Diabetes, Heart Disease, and Parkinson's Disease are part of the project using machine learning methods.

- The Heart Disease model predicts heart disease based on features such as age, blood pressure, cholesterol, and other medical factors.
- The Diabetes model predicts if the individual is diabetic or not by considering factors like glucose, insulin, BMI, and skin thickness.

- The Parkinson's Disease model makes predictions based on clinical data points such as vocal characteristics and motor function measures for the likelihood of Parkinson's.

Scikit-Learn and Pandas libraries are used to train the models, and accuracy scores are used to measure the predictions.

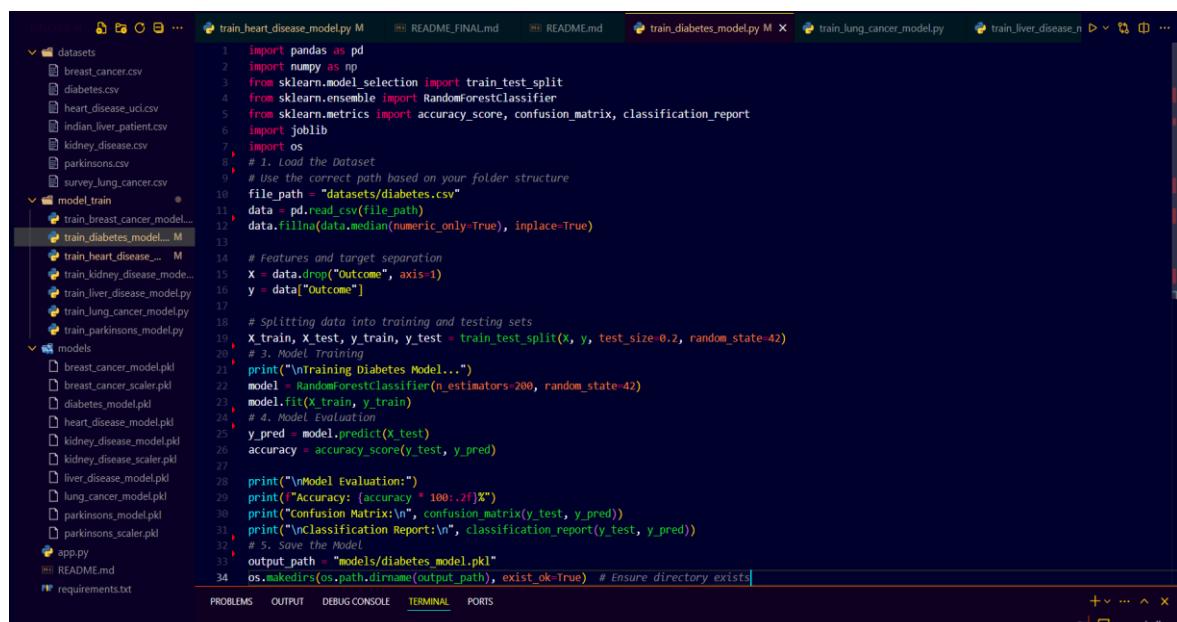
Following are some code segments from the model development process:



```

1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
6 from sklearn.preprocessing import LabelEncoder
7 import joblib
8 # 1. Load the dataset
9 file_path = "datasets/heart_disease_uci.csv"
10 data = pd.read_csv(file_path)
11 data.fillna(data.median(numeric_only=True), inplace=True)
12
13 # Encode categorical variables
14 categorical_cols = ['sex', 'cp', 'fbs', 'restecg', 'exang', 'slope', 'thal']
15 label_encoders = {}
16
17 for col in categorical_cols:
18     le = LabelEncoder()
19     data[col] = le.fit_transform(data[col])
20     label_encoders[col] = le
21
22 # Binarize the target column: 0 (no disease), 1+ (disease + 1)
23 data['target'] = (data['num'] > 0).astype(int)
24
25 # Separating features and labels
26 X = data.drop(['id', 'dataset', 'num'], axis=1) # Remove unnecessary columns
27 y = data['target']
28
29 # Splitting data into training and testing sets
30 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
31 print("\nTraining Heart Disease Model...")
32 model = RandomForestClassifier(n_estimators=200, random_state=42)
33
34 # 4. Model Evaluation
35 y_pred = model.predict(X_test)

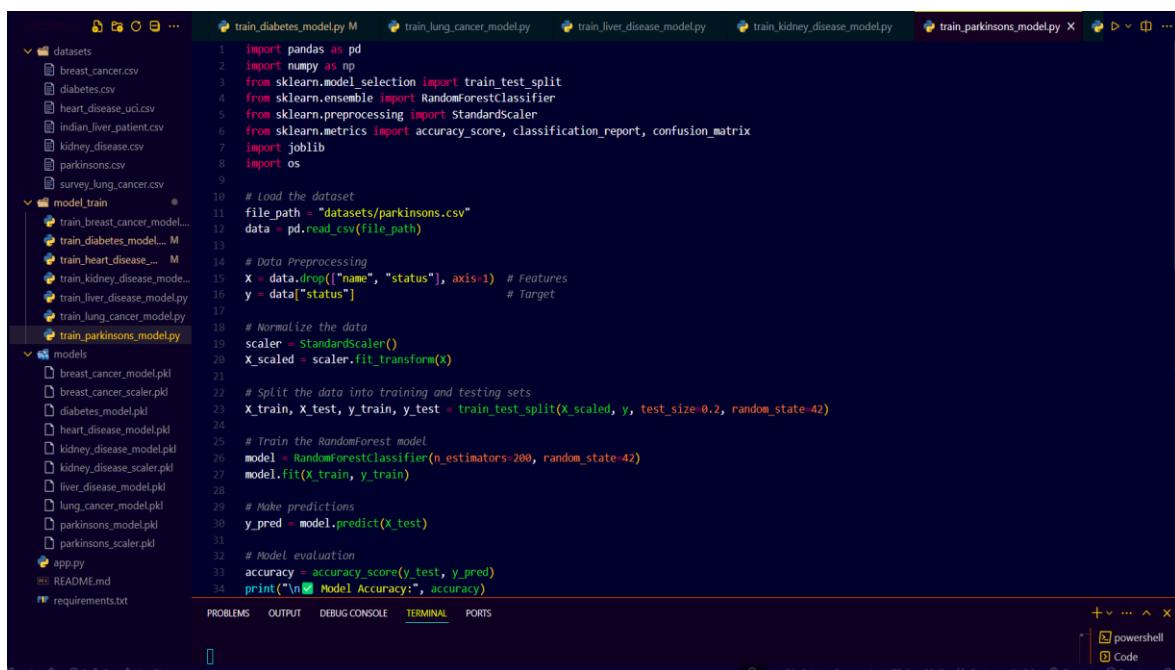
```



```

1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
6 import joblib
7 import os
8 # 1. Load the Dataset
9 # Use the correct path based on your folder structure
10 file_path = "datasets/diabetes.csv"
11 data = pd.read_csv(file_path)
12 data.fillna(data.median(numeric_only=True), inplace=True)
13
14 # Features and target separation
15 X = data.drop("Outcome", axis=1)
16 y = data["Outcome"]
17
18 # Splitting data into training and testing sets
19 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
20
21 print("\nTraining Diabetes Model...")
22 model = RandomForestClassifier(n_estimators=200, random_state=42)
23 model.fit(X_train, y_train)
24
25 y_pred = model.predict(X_test)
26 accuracy = accuracy_score(y_test, y_pred)
27
28 print("\nModel Evaluation:")
29 print("Accuracy: " + str(accuracy * 100))
30 print("Confusion Matrix:\n" + str(confusion_matrix(y_test, y_pred)))
31 print("\nClassification Report:\n" + str(classification_report(y_test, y_pred)))
32
33 output_path = "models/diabetes_model.pkl"
34 os.makedirs(os.path.dirname(output_path), exist_ok=True) # Ensure directory exists

```



```

train_diabetes_model.py M train_lung_cancer_model.py train_liver_disease_model.py train_kidney_disease_model.py train_parkinsons_model.py X train_parkinsons_model.py ...
datasets
breast_cancer.csv
diabetes.csv
heart_disease_uci.csv
indian_liver_patient.csv
kidney_disease.csv
parkinsons.csv
survey_lung_cancer.csv
model_train
train_breast_cancer_model...
train_diabetes_model...
train_heart_disease...
train_kidney_disease...
train_liver_disease_model.py
train_lung_cancer_model.py
train_parkinsons_model.py
models
breast_cancer_model.pkl
breast_cancer_scaler.pkl
diabetes_model.pkl
heart_disease_model.pkl
kidney_disease_model.pkl
kidney_disease_scaler.pkl
liver_disease_model.pkl
lung_cancer_model.pkl
parkinsons_model.pkl
parkinsons_scaler.pkl
app.py
README.md
requirements.txt

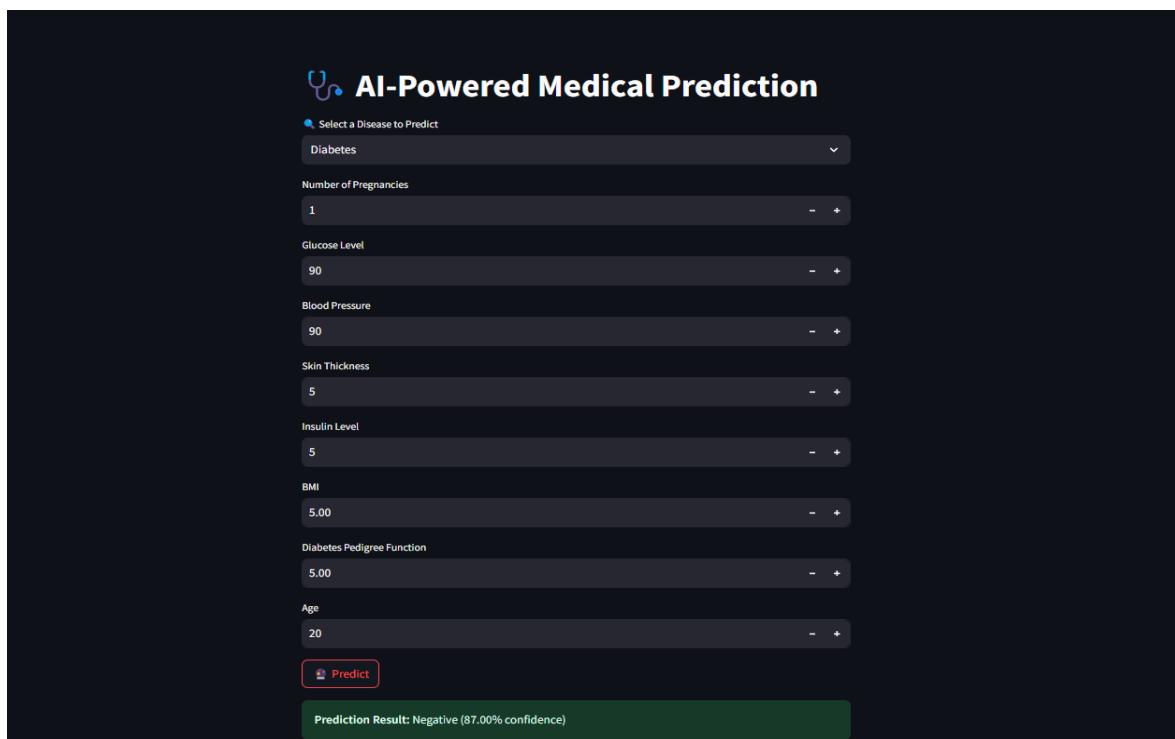
```

TERMINAL

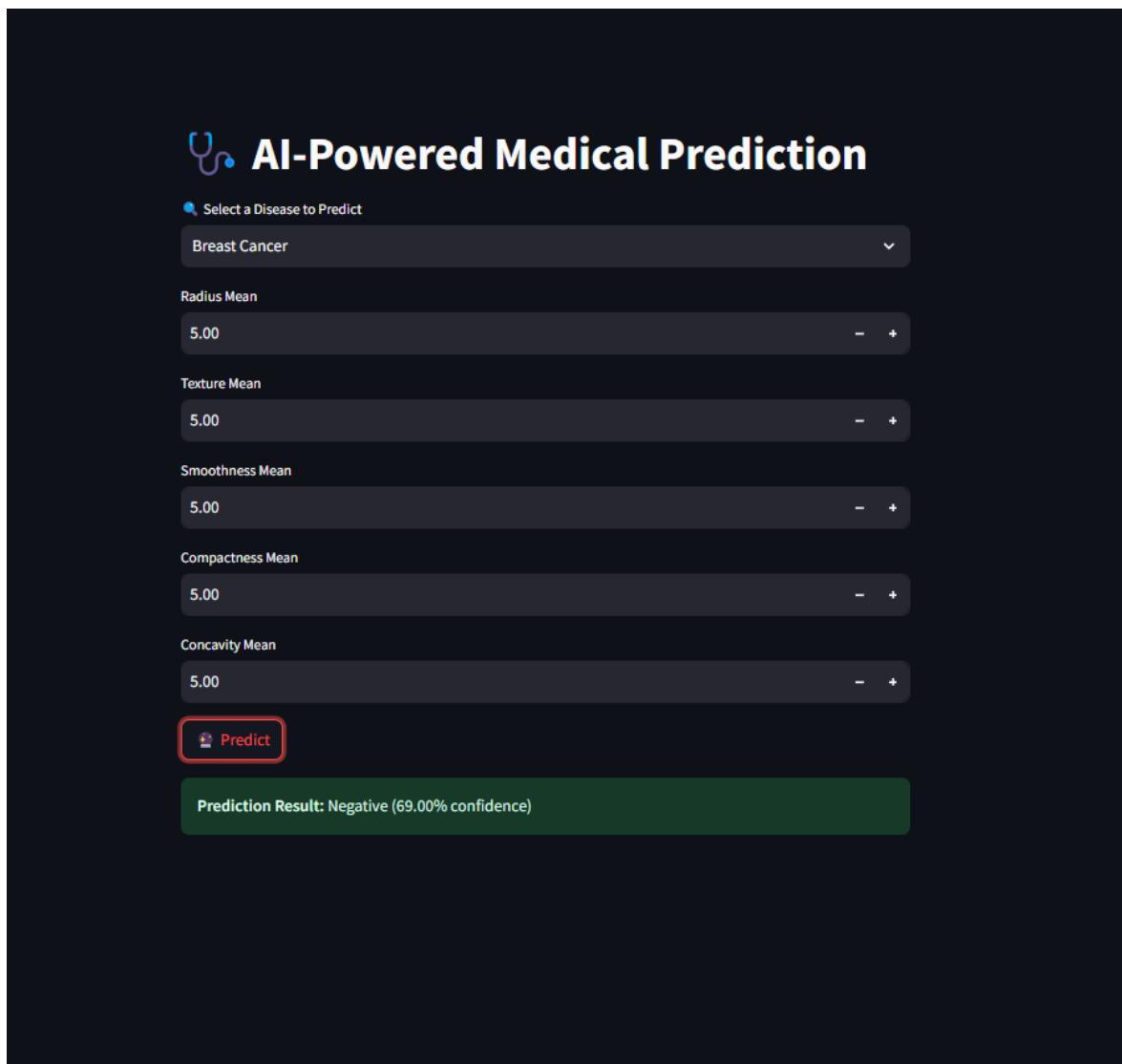
4.2 Web App Development Using Streamlit

A Streamlit-based web application is created using app.py, allowing users to enter health parameters and receive predictions

Diabetes Prediction :



Breast Cancer :



The screenshot shows a web-based medical prediction tool. At the top, there's a logo featuring a stylized blue 'U' and 'S' followed by the text 'AI-Powered Medical Prediction'. Below this, a dropdown menu is set to 'Breast Cancer'. The interface then lists six input fields, each with a numerical value of '5.00' and +/- adjustment buttons: 'Radius Mean', 'Texture Mean', 'Smoothness Mean', 'Compactness Mean', and 'Concavity Mean'. A red-bordered 'Predict' button is located below these inputs. At the bottom, a green bar displays the result: 'Prediction Result: Negative (69.00% confidence)'.

Heart Disease :

就医助手

AI-Powered Medical Prediction

Select a Disease to Predict

Heart Disease

Age: 20

Sex: Male

Chest Pain Type: 1

Resting Blood Pressure: 50

Cholesterol Level: 100

 Predict

Parkinson's Disease :

就医助手

AI-Powered Medical Prediction

Select a Disease to Predict

Parkinson's Disease

MDVP:Fo(Hz): 5.00

MDVP:Fhi(Hz): 5.00

MDVP:Flo(Hz): 5.00

MDVP:Jitter(%): 5.00



4.3 GitHub Link for Code

The complete source code for the project, including model files, datasets, and the Streamlit web app, is available on GitHub.

👉 GitHub Repository:

<https://github.com/aaniket21/AI-Powered-Medical-Prediction>

CHAPTER 5

Discussion and Conclusion

5.1 Future Work

To improve the functionality and performance of the existing system, the following future enhancements are suggested:

- Improving prediction accuracy with deep learning:
 - Applying more sophisticated algorithms like Deep Neural Networks (DNNs), Recurrent Neural Networks (RNNs), and Convolutional Neural Networks (CNNs).
 - Applying these deep learning models can greatly enhance the prediction accuracy and identify more intricate patterns in the health data.
- Incorporating real-time health data for dynamic analysis:
 - Adding the capability for real-time patient data collection via IoT devices, wearable sensors, or health APIs.
 - This will make the system more practical and useful for healthcare professionals by allowing dynamic and real-time predictions.
- Increasing model capabilities to predict other diseases:
 - Adding models for more diseases such as Kidney Disease, Liver Disease, Alzheimer's, and Lung Cancer.
 - This will enhance the usability and versatility of the system for wider medical use.

5.2 Conclusion

This project effectively proves the use of machine learning methods to forecast disease outbreaks. The system makes accurate predictions for Diabetes, Heart Disease, and Parkinson's Disease by applying models learned on medical datasets.

Main Takeaways:

- The Streamlit web application has a friendly user interface to make it easily interactive and hence accessible to medical professionals and researchers.
- The machine learning models created through this project provide useful insights through the detection of patterns in patient data, enabling healthcare professionals to make informed, data-driven decisions.
- This solution has the potential to be scaled to enable real-time diagnosis and extend its ability to include other diseases.

In general, the project demonstrates the potential of machine learning to revolutionize healthcare through early disease detection and facilitating good medical decision-making.

REFERENCES

- Shin H, Cho S. Neighborhood property-based pattern selection for support vector machines. *Neural Comput.* 2007;19:816–855. doi: 10.1162/neco.2007.19.3.816.
- Yang, W., & Lipsitch, M. (2013). "Real-time tracking and forecasting of infectious diseases." *Nature Communications*, 4(1), 2798
- Rajkomar, A., Dean, J., & Kohane, I. (2019). *Machinelearningin medicine*. **New England Journal of Medicine*, 380*(14),1347-1358
- Anderson, R. M., & May, R. M. (1991). *Infectious diseasesofhumans: Dynamics and control*. Oxford University Press.