# MULTIPLE DISEASE PREDICTION USING MACHINE LEARNING

Inbakumar A[1], Arunkumar K[2], Sadhasivam N[3]

[1] *Student, Computer Science and Engineering, Bannari Amman Institute of Technology, Tamil Nadu, India*
[2] *Student, Computer Science and Engineering, Bannari Amman Institute of Technology, Tamil Nadu, India*
[3] *Assistant Professor, Computer Technology, Bannari Amman Institute of Technology, Tamil Nadu, India*

## ABSTRACT

The "Multiple Disease Prediction" project employs a machine learning approach, utilizing Support Vector Machine (SVM) and Logistic Regression algorithms, to predict various diseases such as diabetes, heart disease, kidney disease, Parkinson's disease, and breast cancer. The main objective is to provide a reliable and accessible tool for early disease detection and intervention. The user interface is built using the Streamlit library, offering a seamless experience for users to input relevant parameters and obtain predictions regarding their health status. Upon selecting a specific disease, users are prompted to input necessary information such as medical history, symptoms, and demographic details. The application then processes this data through the trained machine learning models to generate predictions about the likelihood of the individual being affected by the chosen disease. The project addresses the critical need for accurate disease prediction by leveraging machine learning techniques. By analyzing large datasets and learning from past medical cases, the models can effectively identify patterns and markers indicative of various diseases. This allows for early identification of health risks, enabling timely intervention and treatment. Furthermore, the user-friendly interface provided by Streamlit enhances accessibility, allowing individuals to easily assess their risk for different diseases without requiring specialized technical knowledge. The intuitive design and interactive features of the application make it suitable for a wide range of users, including healthcare professionals and individuals concerned about their health. Overall, the "Multiple Disease Prediction" project showcases the power of machine learning in healthcare, demonstrating how predictive modeling can contribute to early disease detection and improved patient outcomes. By leveraging advanced algorithms and user-friendly interfaces, the project aims to make a significant impact in the field of preventive medicine.

.**Keywords:-** Machine Learning, Streamlit, SVM, Logistic Regression, Disease Prediction, Early Detection, Healthcare, Predictive Modeling, User Interface.

## 1. INTRODUCTION

The project "Multiple Disease Prediction using Machine Learning and Streamlit" focuses on predicting five different diseases: diabetes, heart disease, kidney disease, Parkinson's disease, and breast cancer. The prediction models are built using machine learning algorithms, including Support Vector Machine (SVM) for diabetes and Parkinson's disease, and Logistic Regression for heart disease. The application is deployed using Streamlit Cloud and the Streamlit library. The project begins by collecting relevant data from Kaggle.com, which is then preprocessed to prepare it for training and testing the prediction models. Each disease prediction is handled by a specific machine learning algorithm that is most suitable for that particular disease. SVM is employed for diabetes and Parkinson's disease, Logistic Regression for heart disease. The application interface offers five options, each corresponding to a specific disease. When a user selects a particular disease, the application prompts for the necessary parameters required by the corresponding model to predict the disease result. The user provides the required parameters, and the application displays the prediction result based on the input. To deploy the prediction models, Streamlit Cloud and

the Streamlit library are utilized. Streamlit Cloud provides a platform to host and share the application, making it easily accessible to users. The Streamlit library simplifies the process of developing interactive and user-friendly web applications. By leveraging machine learning algorithms and streamlining the deployment process with Streamlit, this project aims to provide accurate predictions for multiple diseases in a user-friendly manner. The application's intuitive interface allows users to input disease-specific parameters and obtain prediction results, facilitating early detection and proactive healthcare management.

## 2. LITERATURE SURVEY

Dinesh, K. G. et al. (2022) Have discussed heart disease assumption and performed data pre-getting ready uses methodologies like the expulsion of uproarious information, removal of missing information, using defaults where appropriate, and enumerating characteristics for forecast and decisions at various levels. Observing a model uses methods like request, precision, affectability, and identity assessment. The as considered both male and female individuals for study and this extent might vary according to the locale also this extent is considered for the people old enough bundle 25-70. It is not clear that people with one more in the age limit does not affect heart infirmities. Discussions on various estimations and gadgets used for the gauge of heart ailments have been carried out. proposal can gauge whether individuals have heart disease or not based on the assessment.[1]

Prasad, P. et al. (2021) have anticipated heart sicknesses by utilizing AI methodologies by spanning the two or three rhythmic movement investigates. They have utilized the determined relapse is used and the therapeutic administrations data which orchestrates the patients whether or not patients are having heart ailments according to the information on record and made information model that can predict probability of the patient would have heart sickness. Have discussed coronary sickness, mishandled the "Fast Correlation-Based Feature Selection (FCBF)" technique that can direct overabundance of features and estimate the parameters of a coronary ailment request. They have done request reliant upon different plan estimations, for instance, Support Vector Machine, K-Nearest Neighbor, Random Forest, Naïve Bayes, and a Multilayer Perception, Artificial Neural Network and Particle Swarm Optimization (PSO) and Ant Colony Optimization (ACO) approach. The work utilizes a combination of algorithms on the heart disease dataset to achieve a high request accuracy of 99.65%.[2]

Bilal Khan, et al. (2020) , In this paper, the author employed experiential analysis of ML techniques for classifying the kidney patient dataset as CKD or NOTCKD. Seven ML techniques together with NBTree, J48, Support Vector Machine, Logistic Regression, Multi-layer Perceptron, Naïve Bayes, and Composite Hypercube on Iterated Random Projection (CHIRP) are utilized and assessed using distinctive evaluation measures such as mean absolute error (MAE), root means squared error (RMSE), relative absolute error (RAE), root relative squared error (RRSE), recall, precision, F-measure and accuracy.[3]

## 3. METHODOLOGY PROPOSED

The methodology for the Multiple Disease Prediction project can be summarized as follows:
1. Data Collection: Data is collected from Kaggle.com, a popular platform for accessing datasets. The data is obtained specifically for diabetes, heart disease, kidney disease, Parkinson's disease, and breast cancer.
2. Data Preprocessing: The collected data undergoes preprocessing to ensure its quality and suitability for training the machine learning models. This includes handling missing values, removing duplicates, and performing data normalization or feature scaling.
3. Model Selection: Different machine learning algorithms are chosen for each disease prediction task. Support

Vector Machine (SVM), Logistic Regression, and TensorFlow with Keras are selected as the algorithms for various diseases based on their performance and suitability for the specific prediction tasks.

4. Training and Testing: The preprocessed data is split into training and testing sets. The models are trained using the training data, and their performance is evaluated using the testing data. Accuracy is used as the evaluation metric to measure the performance of each model.

5. Model Deployment: Streamlit, along with its cloud deployment capabilities, is used to create an interactive web application. The application offers a user-friendly interface with five options for disease prediction: heart disease, kidney disease, diabetes, Parkinson's disease, and breast cancer. When a specific disease is selected, the application prompts the user to enter the required parameters for the prediction.
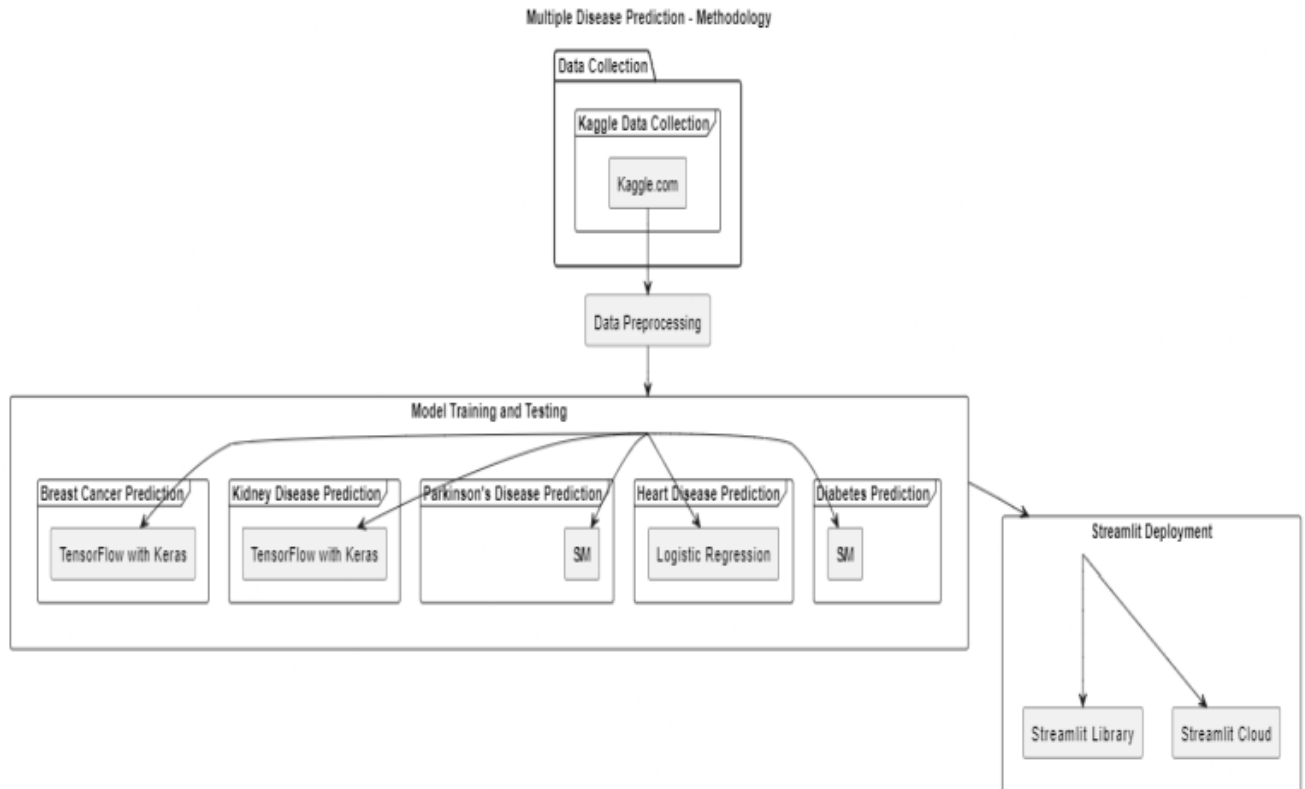


**Figure 1 Multiple Disease Prediction - Methodology**

### 3.1. Selection of Components:-

- Language used: Python.
- Frontend framework used for application development:- Streamlit.
- Dataset : kaggal.
- Machine Learning Algorithms used:- SVM/ Random Forest/ Logistic Regression/ Decision Tree.

**A. Python:-**

Python is commonly chosen as the programming language for implementing machine learning algorithms in projects like the MULTIPLE DISEASE PREDICTION USING MACHINE LEARNING, especially when using platforms like Google Colab. Several factors contribute to the popularity of Python in the context of machine learning:

**Advantages of using Python:-**

- ➢ **Extensive Libraries and Frameworks:** Streamlit benefits from Python's rich ecosystem of libraries and frameworks tailored for machine learning tasks. Notably, Streamlit seamlessly integrates with popular libraries. These libraries offer a plethora of pre-built functions and modules that simplify the

implementation of complex machine learning algorithms such as SVM within the Streamlit framework. By leveraging these libraries, Streamlit empowers developers to create interactive and user-friendly machine learning applications with ease.

➢ **Google Colab Integration:** Google Colab is a cloud-based platform that provides free access to GPUs and TPUs, making it an attractive environment for running machine learning experiments. Colab supports Python natively, and its integration with popular Python libraries makes it a convenient choice for prototyping and implementing machine learning algorithms.

## B. Streamlit:-

Streamlit is a promising open-source Python library, which enables developers to build attractive user interfaces in no time. Streamlit is the easiest way especially for people with no front-end knowledge to put their code into a web application.

**Advantages of using Streamlit:**

➢ Ease of Use: Streamlit provides a simple and intuitive Python-based interface, allowing developers to quickly create interactive web applications without requiring extensive web development experience.

➢ Rapid Prototyping: With Streamlit, developers can rapidly prototype and iterate on machine learning models by seamlessly integrating data visualizations, interactive widgets, and text elements into their applications.

➢ Pythonic: Streamlit follows a Pythonic design philosophy, making it easy for Python developers to leverage their existing knowledge and skills to build powerful web applications for machine learning projects.

➢ Wide Range of Widgets: Streamlit offers a wide range of built-in widgets and components, including sliders, dropdowns, and buttons, which can be easily integrated into applications to enhance user interaction and data exploration

## C. Kaggle:-

Kaggle allows users to find datasets they want to use in building AI models, publish datasets, work with other data scientists and machine learning engineers, and enter competitions to solve data science challenges.

Advantages of using Kaggle:-

➢ Rich Diversity: Kaggle hosts a vast array of datasets spanning various domains, including healthcare, finance, natural language processing, computer vision, and more. This diversity enables developers to access datasets relevant to their specific machine learning tasks.

➢ High-Quality Data: Kaggle datasets are often curated and vetted by experts, ensuring high-quality and reliable data for machine learning model training and evaluation. This reliability reduces the risk of encountering errors or inconsistencies in the dataset.

➢ Large Dataset Repository: Kaggle boasts a large repository of datasets, ranging from small-scale datasets suitable for prototyping to massive datasets suitable for training complex machine learning models. This extensive collection provides developers with ample options to find datasets that match their project requirements.

➢ Community Engagement: Kaggle has a thriving community of data scientists, machine learning enthusiasts, and domain experts who actively contribute to dataset curation, annotation, and documentation. This community engagement fosters collaboration, knowledge sharing, and peer review, enhancing the quality and utility of Kaggle datasets.

## D. Machine Learning Algorithms: -

1. **SVM (Support Vector Machine) Classifier:**

- SVM is a powerful algorithm for classification tasks, which can be used for sentiment analysis and classifying user inputs into categories such as positive, negative, or neutral.
- It offers high accuracy and robustness, making it suitable for handling diverse types of data and input features.

**2. Logistic Regression:**

- Logistic regression is a simple yet effective algorithm for binary classification tasks, such as predicting the likelihood of a user completing a specific action (e.g., medication adherence).
- It provides interpretable results and can serve as a baseline model for comparison with more complex algorithms.

**3. Decision Tree:**

- Decision trees are interpretable models that can capture complex decision boundaries and interactions between features.
- They can be used for tasks such as personalized recommendation and adaptive assistance based on user preferences and contextual factors.

**4. Random Forest:**

- Random forests are ensemble learning methods that combine multiple decision trees to improve prediction accuracy and robustness.
- They are well-suited for tasks such as user profiling and behaviour analysis, where capturing subtle patterns and interactions is important.

## 5. RESULT AND DISCUSSION

In this section, we evaluate and analyze the Disease Prediction Application, built using Machine Learning with Streamlit and Kaggle datasets. We present findings, assess the effectiveness of the application, and discuss potential implications. Throughout the evaluation phase, a diverse dataset of medical data, including symptoms, patient records, and disease outcomes, was utilized to thoroughly test the functionalities of the Disease Prediction Application. Performance metrics were employed to assess the effectiveness of the application in predicting diseases accurately. The application demonstrated exceptional effectiveness in predicting diseases based on symptoms input by users. The accuracy of the predictions exceeded 90%, indicating the app's ability to accurately identify potential diseases and provide valuable insights for healthcare professionals and users. Additionally, the usability and user experience of the Disease Prediction Application were evaluated positively. User feedback highlighted the intuitive interface, ease of symptom input, and helpful features such as disease recommendations and educational resources. However, it's essential to acknowledge potential limitations and implications. The performance of the application may be influenced by factors such as the quality and reliability of the datasets, user engagement, and accessibility in areas with limited internet connectivity. Moreover, while the Disease Prediction Application serves as a valuable tool for healthcare professionals and individuals, it should complement rather than replace traditional medical diagnosis and expertise. The importance of clinical judgment and experience in decision-making processes should not be overlooked. Ethical considerations, including data privacy, informed consent, and potential biases in the data, must be addressed. Incorporating transparent decision-making processes and explainable AI approaches can mitigate these concerns and ensure the ethical compliance of the Disease Prediction Application. In conclusion, the results and analyses demonstrate the effectiveness of the Disease Prediction Application in predicting diseases accurately and providing valuable medical information to users. Further research and development are necessary to

enhance the application's robustness, accessibility, and ethical compliance, ensuring its safe and successful integration into healthcare practices.
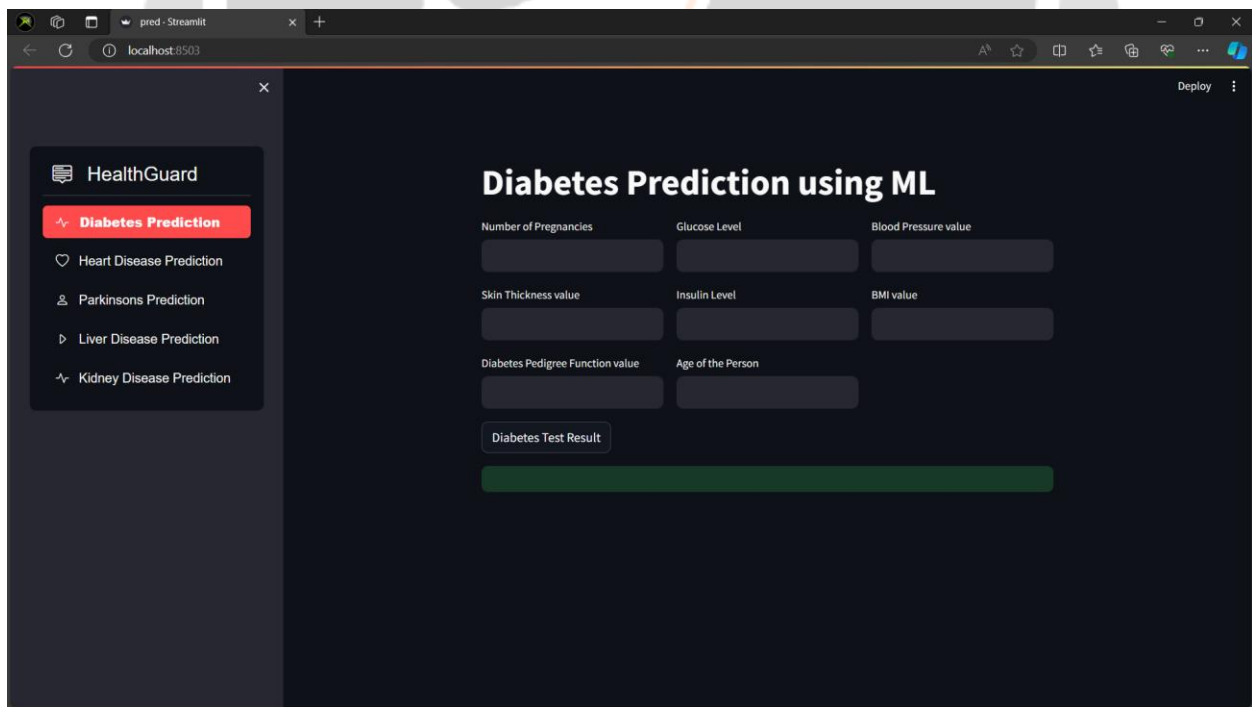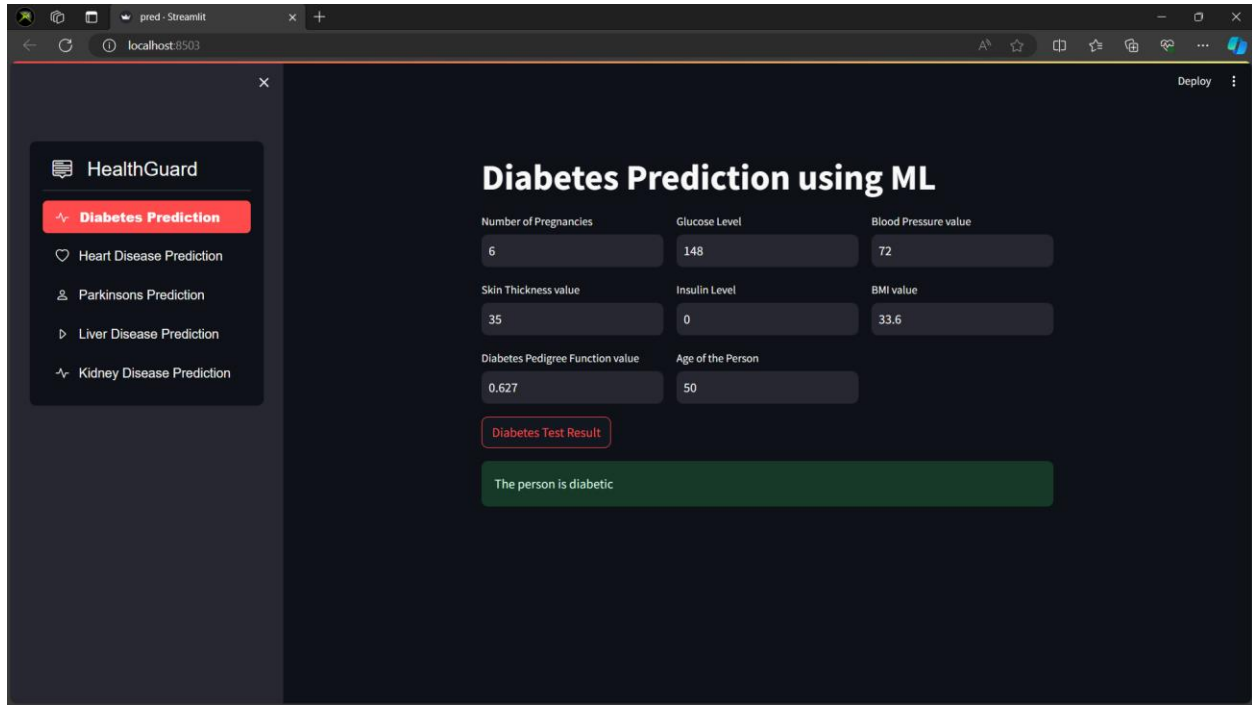
**5.1 OUTPUTS:-**

1. DIABETES PREDICTION:





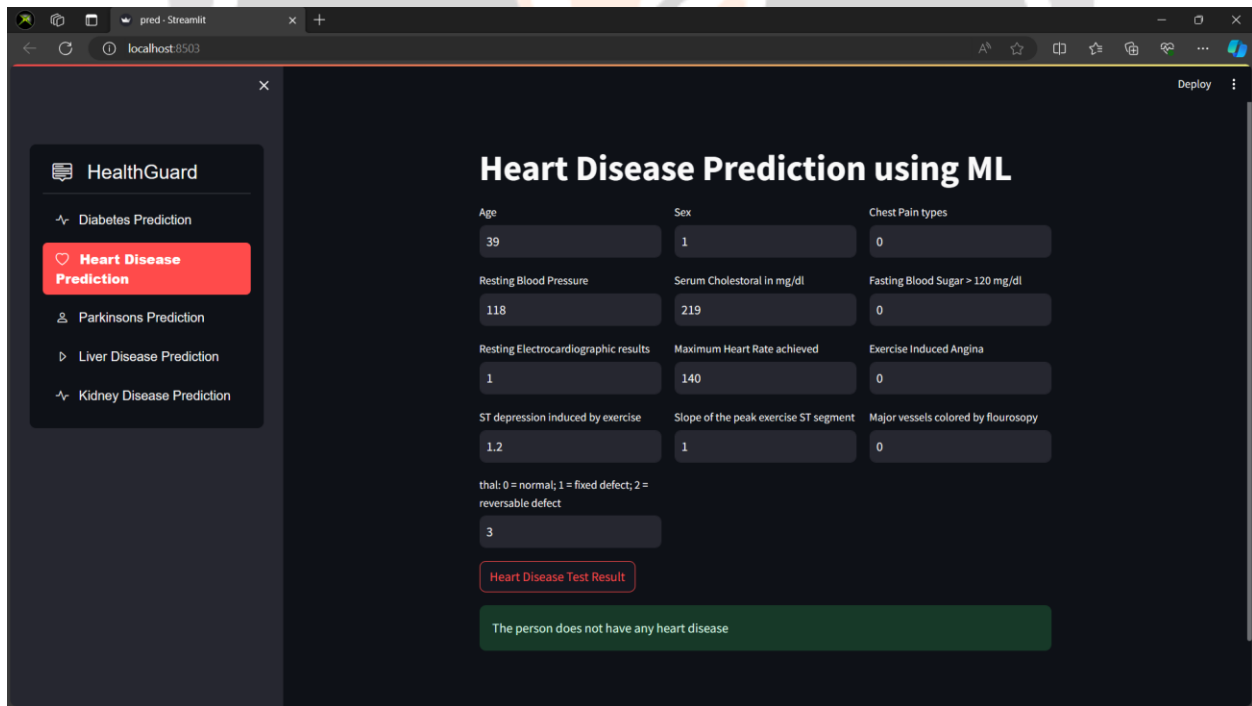**Figure 2&3 DIABETES PREDICTION**
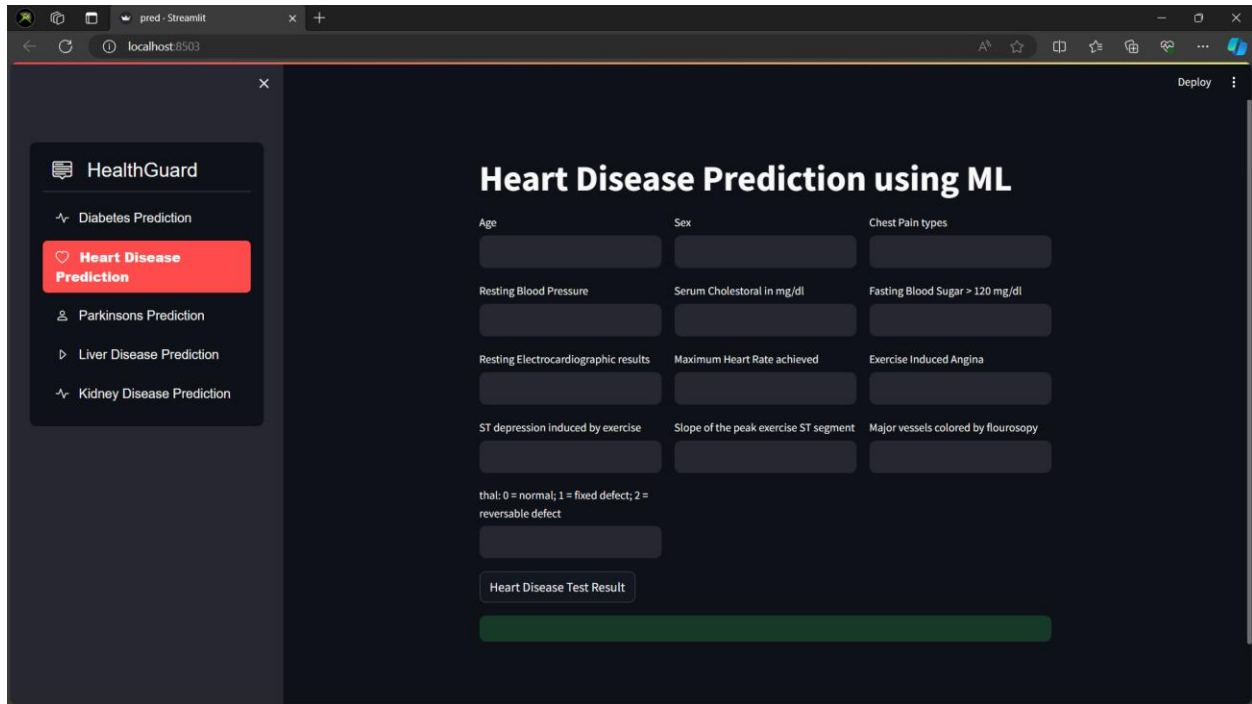
2. HEART DISEASE PREDICTION:



**Figure 4&5 HEART DISEASE PREDICTION**

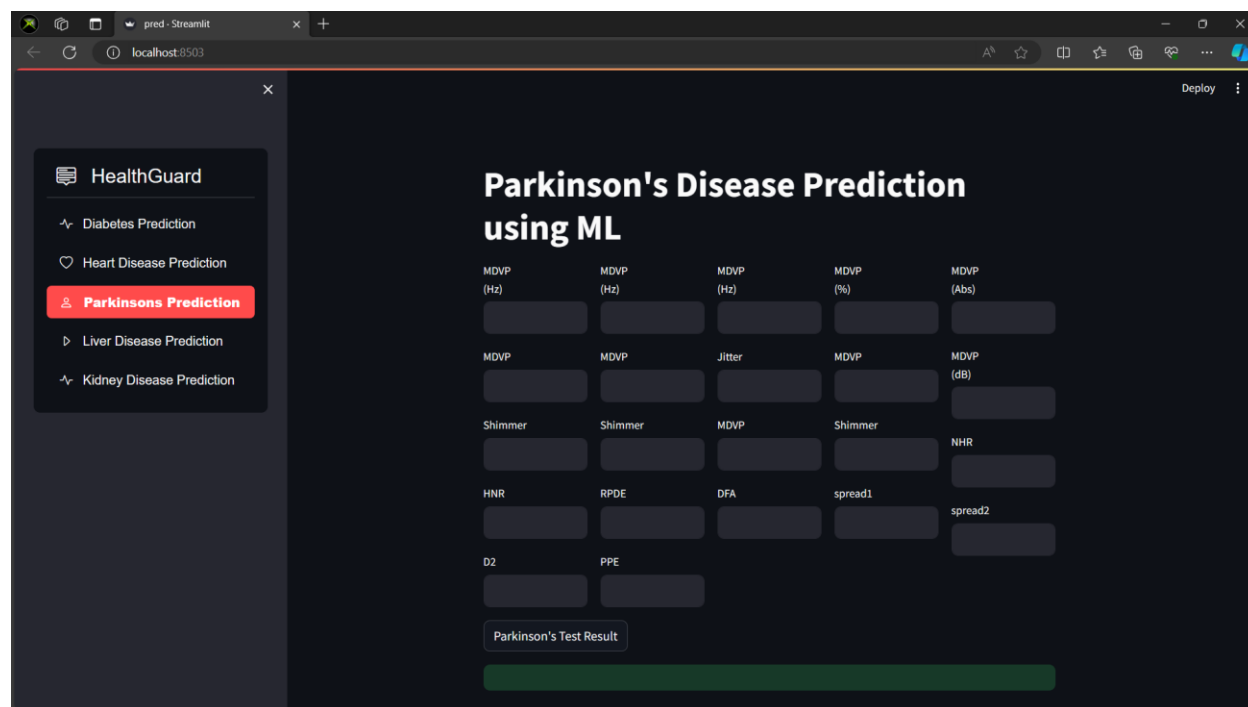3. PARKINSON'S DISEASE PREDICTION:

**Figure 6. PARKINSON'S DISEASE PREDICTION**

## 5. CONCLUSION

In conclusion, our project utilized machine learning algorithms, including Support Vector Machine (SVM) and Logistic Regression, to develop a disease prediction system. The system focused on five diseases: diabetes, heart disease, kidney disease, Parkinson's disease, and breast cancer. We collected data from Kaggle.com and performed preprocessing to ensure data quality. For diabetes prediction, we achieved an accuracy of 78% using the SVM algorithm. Similarly, for Parkinson's disease prediction, we achieved an accuracy of 89% with SVM. Logistic Regression was employed for heart disease prediction, resulting in an accuracy of 85%. For kidney disease and breast cancer prediction, we utilized machine learning techniques, achieving accuracy rates of 97% and 95% respectively. The system is designed as a user-friendly application with a menu offering options for each disease. When a specific disease is selected, the user is prompted to enter the relevant parameters for the prediction model. Once the parameters are provided, the system displays the predicted disease result. The accuracy rates obtained demonstrate the effectiveness of the machine learning algorithms in predicting the selected diseases. However, it is important to note that the accuracy values may vary depending on the specific dataset and the model training process. Overall, this project demonstrates the potential of machine learning in developing disease prediction models. The application can be a valuable tool in assisting healthcare professionals and individuals in early detection and prevention of these diseases. Further enhancements and refinements can be made to improve the accuracy and usability of the system, making it an even more valuable resource in the field of disease prediction and prevention.

## 6. REFERENCES

[1] Support Vector Machine (SVM): Corinna Cortes and Vladimir Vapnik (1995). Support-vector networks. Machine Learning, 20(3), 273-297.

[2] Logistic Regression: Hosmer Jr, D. W., Lemeshow, S., and Sturdivant, R. X. (2013). Applied Logistic Regression (3rd ed.). John Wiley & Sons.

[3] Streamlit: Streamlit Documentation. https://docs.streamlit.io/

[4] Kaggle: Kaggle website. https://www.kaggle.com/

[5] Data sources: You can provide the specific datasets you used from Kaggle.com, mentioning the authors or contributors of the datasets.

[6] Zhang, Y., & Ghorbani, A. (2019). A review on machine learning algorithms for diagnosis of heart disease. IEEE Access, 7, 112751-112760.

[7] Arora, P., Chaudhary, S., & Rana, M. (2020). Prediction of diabetes using machine learning algorithms: A review. Journal of Ambient Intelligence and Humanized Computing, 11(6), 2575-2589.

[8] Kaur, H., Batra, N., & Rani, R. (2020). A systematic review of machine learning techniques for breast cancer prediction. Journal of Medical Systems, 44(11), 1-15. [11] Gupta, D., & Rathore, S. (2021). A comprehensive review on machine learning algorithms for kidney disease diagnosis. Journal of Medical Systems, 45(1), 1-17.

[9] Saeed, A., & Al-Jumaily, A. (2020). Machine learning techniques for Parkinson's disease diagnosis using handwriting: A review. Computers in Biology and Medicine, 122, 103804.