

# Voice Assistant for Disease Diagnosis Using Machine Learning and Natural Language Processing

Dr Swati Sharma  
*Dept. of Computer Science  
Engineering  
Presidency University  
Bengaluru, India*  
swati.sharma@presidencyuniversity.in

Thanusha M  
20191CCE0076  
*Dept. of Computer Science  
Engineering  
Presidency University  
Bengaluru, India*  
201910100204@presidencyuniversity.in

Smitha Reddy S  
20191CCE0061  
*Dept. of Computer Science  
Engineering  
Presidency University  
Bengaluru, India*  
201910100730@presidencyuniversity.in

Sowhardh C K  
20191CCE0065  
*Dept. of Computer Science  
Engineering  
Presidency University  
Bengaluru, India*  
201910100737@presidencyuniversity.in

Shilpa N  
20191CCE0058  
*Dept. of Computer Science  
Engineering  
Presidency University  
Bengaluru, India*  
201910100306@presidencyuniversity.in

**Abstract—** The use of voice assistants in healthcare has become increasingly popular due to their ability to provide remote and personalized care. The proposed idea is to develop a voice assistant model to predict acute diseases using the Random Forest algorithm and make recommendations on treatment and diet plans for the user. The system utilizes the Natural Language Processing (NLP) techniques to record and convert data from Speech to Text and Text to Speech. The Random Forest algorithm is used for feature selection and prediction of diseases based on the symptoms. The openAI library gives access to the openAI Application Programming Interface (API) which is used for information retrieval on treatments and diets. The proposed model can be applied in the development of future healthcare systems that leverage voice assistant technology for improved disease detection and diagnosis. The study highlights the potential of voice assistants in remote areas, making healthcare more accessible and efficient for patients.

**Keywords—** Artificial Intelligence, Machine Learning, Prediction, Natural Language Processing, Voice Assistant, Random Forest

## 1. INTRODUCTION

Artificial Intelligence (AI) is making significant changes in healthcare by offering novel ways to gather and analyze patient information, improve medical decision-making, and personalize medical treatments. One of the most promising AI applications in healthcare is the use of voice assistants, which are powered by machine learning algorithms and natural language processing (NLP), to improve healthcare delivery and patient outcomes. These voice assistants allow patients to communicate their symptoms, obtain medical information, and receive personalized medical advice through voice-based communication.

The use of voice assistants with machine learning algorithms in healthcare can enhance the accessibility, convenience, and efficiency of medical care. These can be particularly beneficial for individuals living in remote areas, where access to medical care may be limited due to

geographical barriers, lack of healthcare facilities, and shortage of medical professionals. Moreover, the technology can alleviate the workload of healthcare professionals, allowing them to focus on critical cases and improving overall healthcare delivery.

Supervised machine learning (ML) models are increasingly being used in disease diagnosis to improve the accuracy of predictions and enhance patient outcomes. The models are trained on labeled datasets and make predictions based on input data. By analyzing large amounts of patient data, these models can identify patterns and risk factors that may not be immediately apparent to human clinicians. Additionally, they can continuously learn and adapt to new data, making them effective tools for disease diagnosis and prediction.

This paper explores the potential of voice assistant systems utilizing ML models and NLP techniques in disease prediction for remote areas healthcare development. It discusses the effectiveness of supervised machine learning models in accurately predicting diseases, leading to more effective interventions and treatments. Furthermore, it highlights the potential of OpenAI in unlocking the full potential of voice assistants in the medical field. Lastly, the paper offers recommendations for future research and development in supervised machine learning models for disease prediction.

## 2. LITRATURE SURVEY

In recent times, healthcare applications have been increasingly adopting machine learning and natural language processing techniques. Among these applications, voice assistants for disease diagnosis have emerged as a promising tool for healthcare providers to remotely monitor patients' health. Several studies have explored the use of machine learning algorithms for disease prediction, various supervised machine learning algorithms for disease prediction are compared in

[2] and it is concluded that Random Forest gave more accuracy, and Rayan Alanazi proposed a machine learning-based approach in [1] for the identification and prediction of chronic diseases using convolution neural networks (CNN) and K-nearest neighbor (KNN). These studies highlight the potential of machine learning in disease prediction, laying the groundwork for the development of voice assistant systems for disease diagnosis.

Furthermore, natural language processing techniques have also been investigated in the development of voice assistant systems. In paper [3] a review is conducted on speech-to-text and text-to-speech recognition systems, emphasizing the importance of natural language processing in enabling voice assistants to accurately interpret human speech. Additionally, the proposed system in [4] is an end-to-end text-to-speech synthesis system that generates speech with human-level quality, which could enhance the user experience of voice assistant systems. Moreover, chatbot systems that utilize natural language processing and artificial intelligence techniques for medical diagnosis have also been explored by paper [5] and paper [6]. Collectively, these studies demonstrate the potential of natural language processing techniques in improving the user experience of voice assistant systems for disease diagnosis.

## 3. PROPOSED SYSTEM

In this section, a detailed description on datasets collection, model development, disease prediction, and voice assistant creation is given. The initial step in constructing a machine learning model is to collect data. Datasets were obtained from Kaggle, a data science platform. After data collection, the data is processed and divided into training and testing datasets. Then the datasets were trained and tested with the machine learning algorithms such as SVM, Naïve Bayes, Decision Trees and Random Forest (RF). And when

compared for accuracy RF is selected. This model is then integrated with the voice assistant program.

Following are the steps involved in creation of Voice Assistant for Disease Diagnosis.

- 3.1. Data Collection.
- 3.2. Data Preprocessing.
- 3.3. Disease Prediction Using Random Forest.
- 3.4. Speech to Text Using SpeechRecognition.
- 3.5. Text to Speech Using Pyttsx3.
- 3.6. Information Retrieval Using OpenAI.

**3.1. Data Collection.** The data collected includes 132 common symptoms(features) mapped to 41

data using the train\_test\_split function from the sklearn library. The split data (symptoms and diseases) is then fitted onto the Random Forest Model to train. Later the model is tested on the test dataset. The illustration of Random Forest algorithm consisting of 3 different decision trees is shown in the Fig 1. Each of the decision tree was trained using a random subset of the training data.

**3.4. Speech to Text Using SpeechRecognition.** This model converts the user's spoken words into text, which allows the voice assistant to understand what the user is saying. It is used to translate the patients' symptoms into a digital format that can be processed and analyzed.

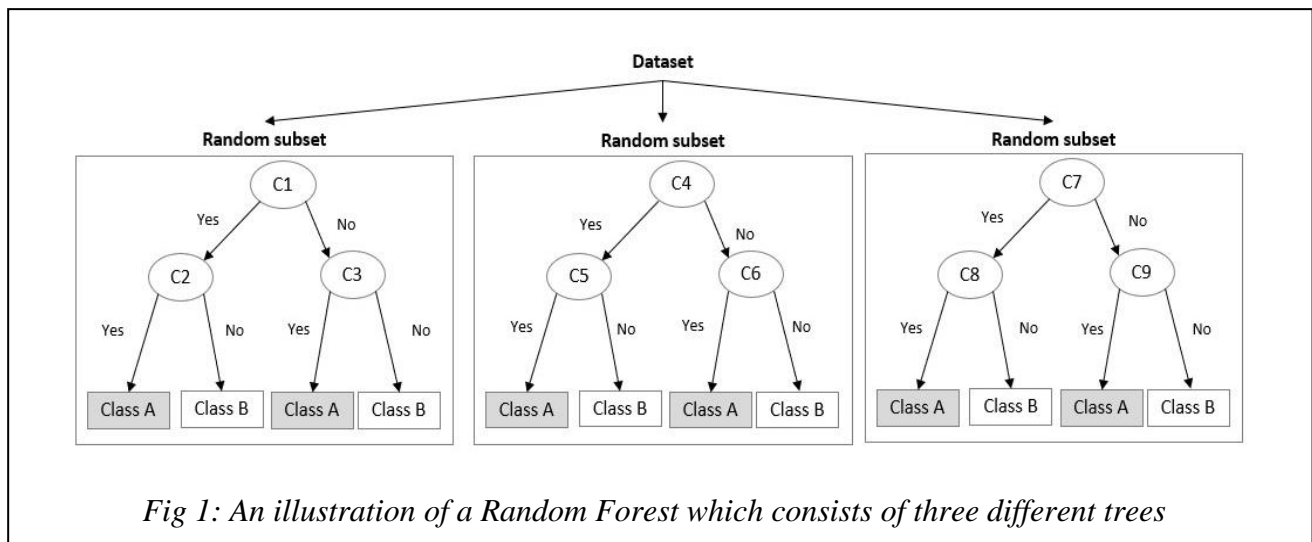


Fig 1: An illustration of a Random Forest which consists of three different trees

unique diseases(target). The dataset excludes personal details of patients such as name, ID, mobile number and so on to prevent privacy.

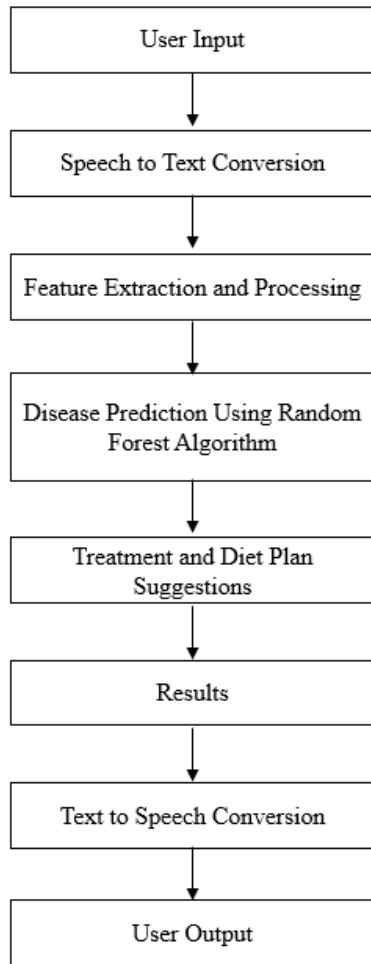
**3.2. Data Preprocessing.** The collected data are preprocessed to check for null values and to drop unnecessary columns. This is done to enhance the quality of the dataset. This step also eliminates underscores, commas, and white spaces. Once the preprocessing is done, it is ready for training and testing.

**3.3. Disease Prediction Using Random Forest.** The proposed system uses Random Forest algorithm to predict the acute diseases. The processed train dataset is split into test and train

**3.5. Text to Speech Using Pyttsx3.** Pyttsx3 is a Python library, used for text-to-speech conversion. When integrated into the voice assistant it can be used to provide spoken responses to users about the results of their diagnosis. And also makes recommendation on treatments and diets.

**3.6. Information retrieval Using OpenAI.** The link to OpenAI GPT model is established with the API key. This further helps in retrieval of information on treatment plans, diet charts and so on from the GPT model.

After all the steps, if the voice assistant's performance matches the desired expectations, then the proposed system is ready for deployment as shown in Fig 2.



*Fig 2: Architecture of proposed Voice Assistant for disease diagnosis system.*

#### 4. EXPERIMENTAL RESULTS

The Experimental results of the proposed Voice Assistant system indicate that it achieved the desired outcomes. the accuracy of the Random Forest algorithm was found to be 92.68%, which is a promising result for predicting acute diseases. The system's speech-to-text and text-to-speech

conversions were accurate, making it easier for patients to communicate with the system. The information retrieved from OpenAI was also useful and accurate, further demonstrating the system's effectiveness in providing patients with reliable and relevant information on treatments and diets.

#### 5. CONCLUSION

In this paper, we proposed a voice assistant system to predict diseases and make recommendations on treatments and diets based on the user's symptoms. It seeks to convert the input audio to machine understandable form using NLP techniques. It makes predictions using the random forest algorithm. Moreover, with the help of OpenAI, the system can provide accurate and reliable suggestions on treatments and diets. The system can prove to be a game-changer in rural areas where health facilities are limited, and people have to travel long distances to get a proper diagnosis. Overall, the proposed system has the potential to bring significant changes to the healthcare industry and rural sectors.

#### 6. FUTURE SCOPE

Future work can involve expanding the dataset to include a wider range of diseases and symptoms, and incorporating other technologies such as image recognition. The study can further be incorporated with neural networks to achieve better understanding of user's input and generate more accurate disease diagnosis. The system only accepts user's input in English language, it can be extended to include other regional and international languages. Future systems could explore the potential use of the system in clinical settings, where it could aid healthcare professionals in making more informed diagnosis and treatment decisions.

## REFERENCES

- [1] Rayan Alanazi, "Identification and Prediction of Chronic Diseases Using Machine Learning Approach", *Hindawi Journal of Healthcare Engineering, Volume 2022*.
- [2] Shahadat Uddin<sup>1</sup>, Arif Khan, Md Ekramul Hossain and Mohammad Ali Moni, "Comparing different supervised machine learning algorithms for disease prediction", *BMC Medical Informatics and Decision Making*, 2019.
- [3] Ayushi Trivedi, Navya Pant, Pinal Shah, Simran Sonik, Supriya Agrawal, "Speech to text and text to speech recognition systems-A review", *IOSR Journal of Computer Engineering*, 2018.
- [4] Xu Tan, Jiawei Chen, Haohe Liu, Jian Cong, Chen Zhang, Yanqing Liu, Xi Wang Yichong Leng, Yuanhao Yi, Lei He, Frank Soong Tao Qin, Sheng Zhao, Tie-Yan Liu, "NaturalSpeech: End-to-End Text to Speech Synthesis with Human-Level Quality", *arXiv preprint arXiv:2205.04421*, 2022.
- [5] Divya S, Indumathi V, Ishwarya S, Priyasankari M, Kalpana Devi S, "A Self-Diagnosis Medical Chatbot Using Artificial Intelligence", *Journal of Web Development and Web Designing Volume 3 Issue 1*, 2018.
- [6] Dr. Meera Gandhi, Vishal Kumar Singh, Vivek Kumar, "IntelliDoctor – AI based Medical Assistant", *Fifth International Conference on Science Technology Engineering and Mathematics*, 2019.
- [7] Nicholas A. I. Omeregbe, Israel O. Ndaman, Sanjay Misra, Olusola O. Abayomi-Alli, Robertas Damasevicius, "Text Messaging Based Medical Diagnosis Using Natural Processing and Fuzzy Logic", *Hindawi Journal of Healthcare Engineering, Volume 2020*.
- [8] Dong Jin Park, Min Woo Park, Homin Lee, Young-Jin Kim, Yeongsic Kim & Young Hoon Park, "Development of machine learning model for diagnostic disease prediction based on laboratory tests", *Nature Portfolio*, 2021.