

Developing a Machine Learning-Based Multiple Disease Prediction System: A Comprehensive Analysis of Risk Factors and Disease Interactions

Emad Naushad¹, Bhavishya Raj², Arpit Nirvan³, Vrinda Sachdeva⁴

Students, B.Tech Computer Science & Engineering^{1,2,3}

Associate Professor of Department of Computer Science & Engineering⁴

I.T.S Engineering College, Greater Noida, India

Abstract: Using predictive modelling, the "Multiple Disease Prediction System" foretells the user's sickness depending on the symptoms are supplied as input to the system. The system evaluates the user's symptoms as input and outputs the likelihood that the disease will occur. The Random Forest Classifier is used for prediction, and Deep Learning Models for Diabetes, Heart Disease and Parkinson's Disease. This method is more accurate and a construction of a web application for prediction system is done.

Keywords: Parkinson's Disease, Predictive Modelling, Diabetes, Heart Disease

I. INTRODUCTION

Anyone who is currently ill must see a doctor, which is both time consuming and costly. Because the sickness cannot be identified, it might also be challenging for the user if it is out of reach of doctors and hospitals. So, if the following treatment can be performed using automated software that saves time and money, it may be better for the patient and make the process move more smoothly. Other Multiple Disease Prediction Systems examine the patient's risk level using data mining approaches. Ailment Predictor is a web-based application that predicts a user's ailment based on their symptoms. For the Disease Prediction system, data sets from several health-related websites were acquired. Using condition Predictor, the customer will be able to predict the possibility of a condition based on the symptoms provided. People are constantly eager to learn new things, especially as the use of the internet expands by the day. When a problem arises, individuals frequently want to search it up on the internet. Hospitals and physicians have less internet connectivity than the general public. People who are plagued with a sickness do not have many options. As a result, people may benefit from this system. Chronic illness is a disease that lasts for an extended period or takes a long time to heal, and many chronic diseases cannot be cured but must be managed daily. India, like all other countries, is undergoing substantial social and economic changes, which is leading to an increase in the prevalence of cardiovascular disease. Many established, developing, and developing countries, including India, are coping with a wide spectrum of chronic diseases, particularly cardiovascular disease, which has major implications for global health, security, and the economy. The world's growing urbanisation and economic progress have resulted in a diverse spectrum of lifestyles. Chronic diseases are now an issue in all countries, affecting one-third of the population in each. Chronic disease treatment is more expensive, and it is challenging on the sick. A vast number of chronic disease datasets are collected and processed in the medical area, and data mining aids in disease early detection. The most expensive diseases to diagnose include heart disease, diabetes, and Parkinson's disease.

Offering the finest quality services to all patients is a huge difficulty in the medical or healthcare industry, and only those who can afford it may profit from it. There is a tremendous amount of healthcare data available that is not being mined in a more efficient and dependable manner to unearth hidden knowledge for effective decision-making. To diagnose chronic diseases early, the suggested system leverages data mining approaches. Machine learning is the process of teaching computers to improve their output based on past data or examples. Machine learning is the study of computer systems that learn from data and experience. The machine learning algorithm has two stages: training and testing. Prediction of a disease based on the patient's symptoms and medical history for decades, machine learning has

been a stumbling barrier. In the medical sector, Machine Learning technology provides a strong venue for rapidly resolving healthcare challenges.

II. RESEARCH OBJECTIVE

There is a need to explore and build a system that would allow end users to forecast chronic diseases without needing to consult a physician or doctor. Identifying various diseases by studying patients' symptoms and employing various Machine Learning Models approaches. There is no standard process for dealing with text and structured data. The proposed approach would consider both organised and unstructured data. Machine Learning can enhance forecast accuracy.

III. LITERATURE REVIEW

Diabetes is a chronic metabolic disorder that affects millions of people in India. According to the International Diabetes Federation, India had 77 million adults (aged 20-79 years) living with diabetes in 2019, and this number is projected to rise to 134 million by 2045 [1]. Literature on diabetes in India indicates that there are various risk factors that can contribute to the development of this disease, including obesity, physical inactivity, and genetics [2]. In addition to these risk factors, the prevalence of diabetes is also associated with socioeconomic status, age, and ethnicity in India [3]. The management of diabetes in India involves lifestyle changes such as diet and exercise, along with medications such as metformin and insulin [4]. In recent years, research in India has also focused on the use of traditional Indian medicines, such as Ayurveda, to manage diabetes [5]. These therapies hold promise for the treatment of diabetes in India, where traditional medicine plays an important role in healthcare. Heart disease is another chronic disease that affects millions of people in India. According to the Indian Heart Association, cardiovascular disease CVD is the leading cause of death in India, responsible for 28% of all deaths [6]. Literature on heart disease in India indicates that there are several risk factors that can contribute to its development, including smoking, high blood pressure, and high cholesterol [7]. In addition to these risk factors, the prevalence of heart disease is also associated with age, sex, and genetics in India [8]. The management of heart disease in India involves lifestyle changes such as diet and exercise, along with medications such as statins and antiplatelet drugs [9]. In recent years, research in India has also focused on the use of traditional Indian medicines, such as Ayurveda, to manage heart disease [10]. These therapies hold promise for the treatment of heart disease in India, where traditional medicine is an integral part of the healthcare system. Parkinson's disease is a progressive neurological disorder that affects millions of people in India. According to a study published in the Journal of Parkinson's Disease, the prevalence of Parkinson's disease in India is 39.4 per 100,000 population [11]. Literature on Parkinson's disease in India indicates that there are several risk factors that can contribute to its development, including genetics, environmental factors, and aging [12]. In addition to these risk factors, the prevalence of Parkinson's disease is also associated with sex, with men being more likely to develop the disease than women in India [13]. The management of Parkinson's disease in India involves medications such as levodopa and dopamine agonists, along with physical therapy and deep brain stimulation [14]. In recent years, research in India has also focused on the use of traditional Indian medicines, such as Ayurveda, to manage Parkinson's disease [15]. These therapies hold promise for the treatment of Parkinson's disease in India, where traditional medicine is an important part of the healthcare system.

This research paper was written by Emad Naushad, Bhavishya Raj, Arpit Nirvan and Vrinda Sachdeva to provide a survey of existing techniques of information discovery in databases using data mining techniques that are used in today's medical research, specifically in Multiple Disease Prediction. Several experiments have been carried out to compare the performance of predictive modelling techniques on the same dataset, and the results show that Decision Tree outperforms, with Bayesian classification having comparable accuracy to Decision Tree in some cases, but other predictive approaches such as SVM, Logistic Regression, and Classification based on Clustering underperform.

A study was conducted to predict heart diseases using the Decision Tree Algorithm, in which the consumer provides data that is compared to a qualified set of values. As a result of this study, patients were able to provide basic information that was compared to data, and heart disease was expected. Also, analysis of the various types of heart-related problems using medical data mining techniques such as association rule mining, grouping, and clustering I. The aim of a decision tree is to show any possible outcome of a decision. To achieve the best result, various rules are

devised. The criteria used in this study were age, sex, smoking, being overweight, drinking alcohol, blood sugar, heart rate, and blood pressure.

IV. PROPOSED SYSTEM

Architecture Diagram

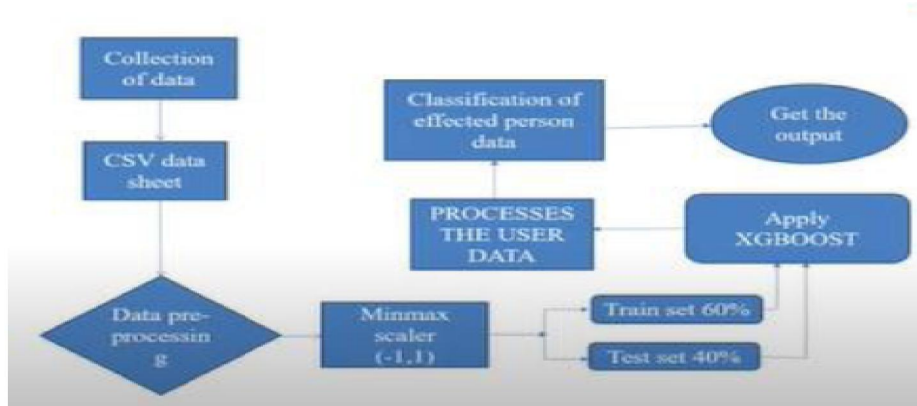


Fig.1. Architecture Diagram

Work Flow Diagram

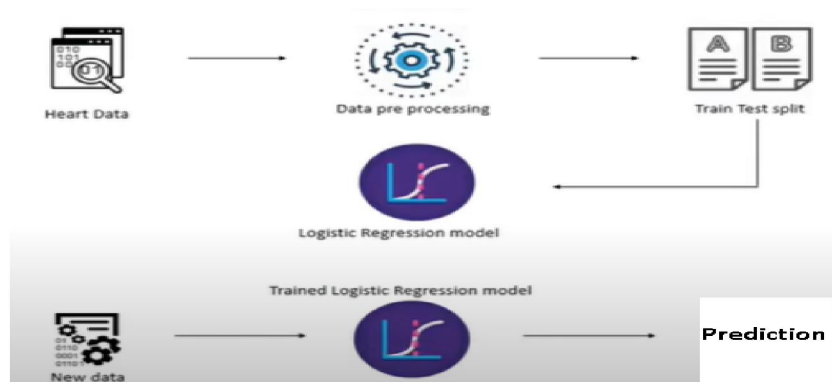
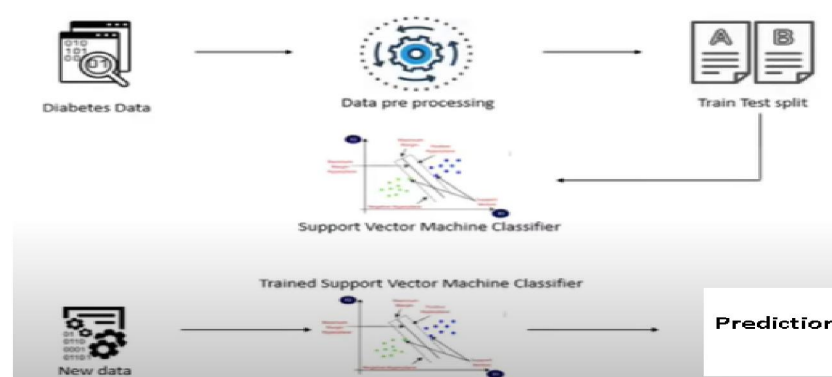


Fig.2. Work Flow Diagram

4.1 Data Collection



Data collection has been done from the internet to identify the disease here the real symptoms of the disease are collected i.e., no dummy values are entered. The symptoms of the disease are collected from different health related websites.

4.2 Data Pre-Processing

Before feeding the data into the Prediction model, following data cleaning and pre-processing steps are performed
Checking null values and filling using forward fill method

4.3 Building Model

- Converting data into different cases
- Standardizing the data using mean and standard deviation
- Splitting the dataset into training and testing sets

Many methods are used to perform data mining. Machine learning is one of the approaches. Random forest Machine learning strategies include grouping, clustering, summarization, and many others. Since classification techniques are used in this project, classification is one of the data mining processes in this phase of categorical data classification. And this step is divided into two phases: training and testing. In the training phase, predetermined data and associated class labels are used for classification. The training stage is often referred to as supervised learning. The preparation and testing phases of the classification process are depicted in the diagram. In the training process, training tuples are used, and in the test data phase, test data tuples are used, and the classification rule's accuracy is calculated. Assume that the classification rule's accuracy on testing data is sufficient for the rule to be used for classification of unmined data.

4.4 Prediction

Prediction using Random Forest: -

Prediction done by Random Forest Model using Streamlit framework model trained by training chronic disease dataset.

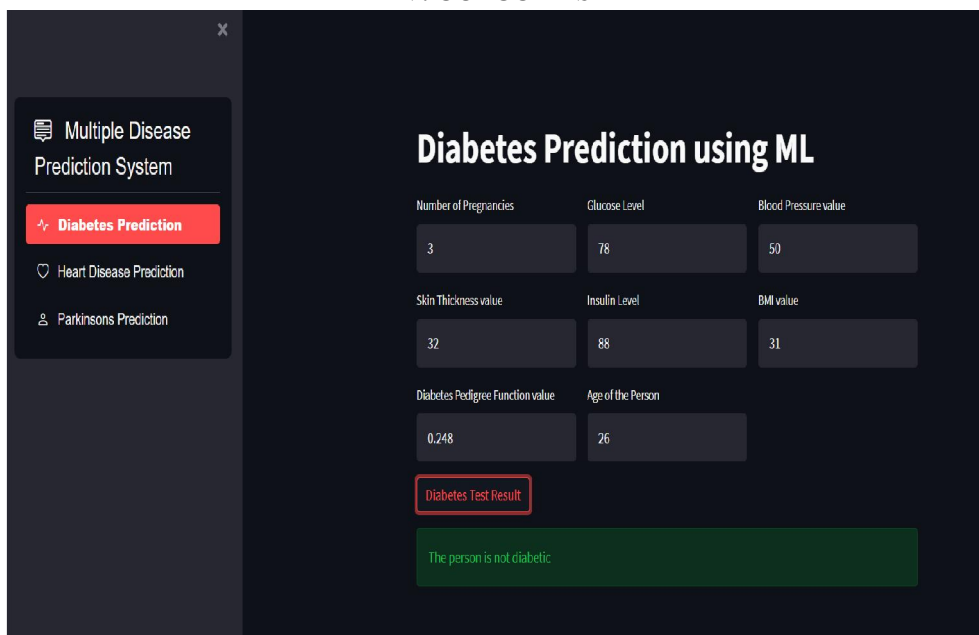
4.5 Algorithm

Logistic regression predicts the output of a categorical dependent variable.

Therefore, the outcome must be a categorical or discrete value. It can be either Yes or No, 0 or 1, true or False, etc. but instead of giving the exact value as 0 and 1, it gives the probabilistic values which lie between 0 and 1.

SVM (Support Vector Machine) algorithm creates the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane.

V. OUTCOMES



Parkinson's Disease Prediction using ML

Multiple Disease Prediction System

- Diabetes Prediction
- Heart Disease Prediction
- Parkinsons Prediction

MDVP:F0(Hz)	MDVP:F1(Hz)	MDVP:F2(Hz)	MDVP:Jitter(%)	MDVP:Jitter(Abs)
119.992	157.302	74.997	0.00784	0.00007
MDVP:RAP	MDVP:PPQ	Jitter:DDP	MDVP:Shimmer	MDVP:Shimmer(dB)
0.0037	0.00554	0.01109	0.04374	0.426
Shimmer:APQ3	Shimmer:APQ5	MDVP:APQ	Shimmer:DDA	NHR
0.02182	0.0313	0.02971	0.06545	0.02211
HNR	RPDE	DFA	spread1	spread2
21.033	0.41478	0.81529	-4.813	0.26648
D2	PPE			
2.30144	0.28465			

Parkinson's Test Result

The person has Parkinson's disease

Heart Disease Prediction using ML

Multiple Disease Prediction System

- Diabetes Prediction
- Heart Disease Prediction
- Parkinsons Prediction

Age	Sex	Chest Pain types
63	1	3
Resting Blood Pressure	Serum Cholesterol in mg/dl	Fasting Blood Sugar > 120 mg/dl
145	233	1
Resting Electrocardiographic results	Maximum Heart Rate achieved	Exercise Induced Angina
0	150	0
ST depression induced by exercise	Slope of the peak exercise ST segment	Major vessels colored by flourosopy
2.3	0	0

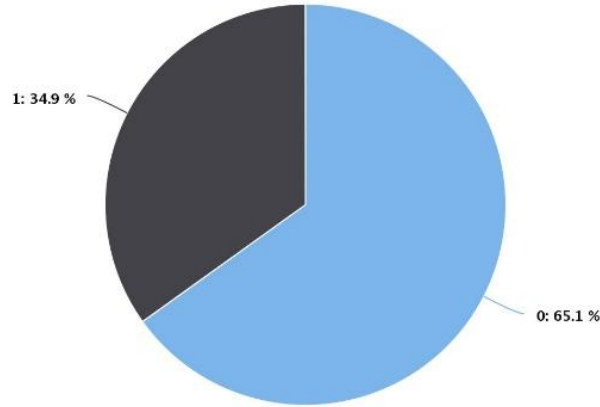
that: 0 = normal; 1 = fixed defect; 2 = reversible defect

Heart Disease Test Result

The person is having heart disease

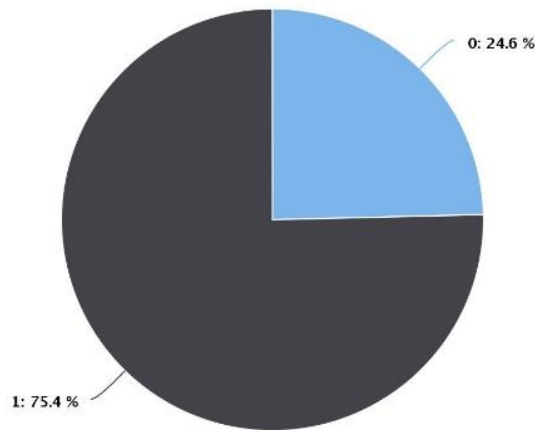
VI. CHARTS

Outcome



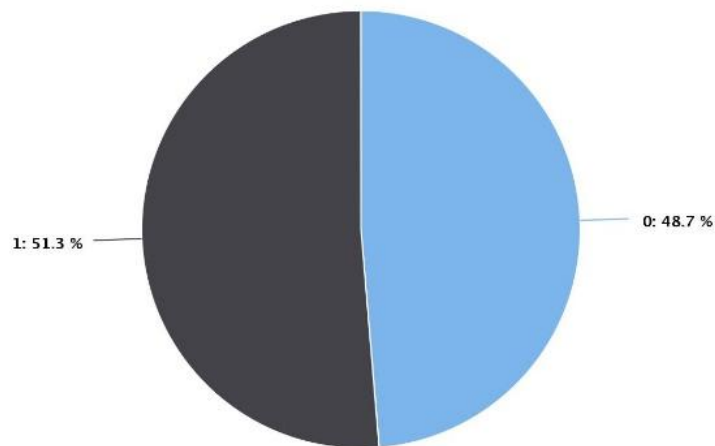
Highcharts.com

status



Highcharts.com

target



Highcharts.com

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target	
2	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1	
3	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1	
4	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1	
5	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1	
6	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1	
7	57	1	0	140	192	0	1	148	0	0.4	1	0	1	1	
8	56	0	1	140	294	0	0	153	0	1.3	1	0	2	1	
9	44	1	1	120	263	0	1	173	0	0	2	0	3	1	
10	52	1	2	172	199	1	1	162	0	0.5	2	0	3	1	
11	57	1	2	150	168	0	1	174	0	1.6	2	0	2	1	
12	54	1	0	140	239	0	1	160	0	1.2	2	0	2	1	
13	48	0	2	130	275	0	1	139	0	0.2	2	0	2	1	
14	49	1	1	130	266	0	1	171	0	0.6	2	0	2	1	
15	64	1	3	110	211	0	0	144	1	1.8	1	0	2	1	
16	58	0	3	150	283	1	0	162	0	1	2	0	2	1	
17	50	0	2	120	219	0	1	158	0	1.6	1	0	2	1	
18	58	0	2	120	340	0	1	172	0	0	2	0	2	1	
19	66	0	3	150	226	0	1	114	0	2.6	0	0	2	1	
20	43	1	0	150	247	0	1	171	0	1.5	2	0	2	1	
21	69	0	3	140	239	0	1	151	0	1.8	2	2	2	1	
22	59	1	0	135	234	0	1	161	0	0.5	1	0	3	1	
23	44	1	2	130	233	0	1	179	1	0.4	2	0	2	1	
24	42	1	0	140	226	0	1	178	0	0	2	0	2	1	
25	61	1	2	150	243	1	1	137	1	1	1	0	2	1	
26	40	1	3	140	199	0	1	178	1	1.4	2	0	3	1	
27	71	0	1	160	302	0	1	162	0	0.4	2	2	2	1	
28	59	1	2	150	212	1	1	157	0	1.6	2	0	2	1	
29	51	1	2	110	175	0	1	123	0	0.6	2	0	2	1	

Heart disease dataset

VIII. CONCLUSION

The proposed work brings diabetes, heart disease, and Parkinson disease under a single platform by deploying the trained models using the Streamlit framework which is a lightweight framework. One classification and one regression algorithms are used for training the models, in which the SVM gave good accuracy values for the disease prediction of diabetes and Parkinson Logistic regression for the disease prediction of heart disease. Its highest accuracy is calculated by picking the highest value obtained from 1 to 21 neighbours. In the future, we can expand this work by adding more diseases that are trained by machine learning models and can include the disease that involves deep learning models.

REFERENCES

- [1]. Agardh E, Allebeck P, Hallqvist J, Moradi T, Sidorchuk A. Type 2 diabetes incidence and socio-economic position: a systematic review and meta-analysis. *Int J Epidemiol.* 2011;40:804–818.
- [2]. Kalia, L.V.; Lang, A.E. Parkinson’s Disease. *Lancet* 2015, 386, 896–912. 4. D. Heisters, "Parkinson's: symptoms treatments and research", vol. 20, no. 9, pp. 548-554, 2011.
- [3]. Zhilbert Tafa, Nerxhivane Pervetica, Bertran Karahoda, “An Intelligent System for Diabetes Prediction”, 4thMediterranean Conference on Embedded Computing MECO – 2015 Budva, Montenegro
- [4]. Mahlknecht, P.; Krismer, F.; Poewe, W.; Seppi, K. Meta-Analysis of Dorsolateral Nigral Hyperintensity on Magnetic Resonance Imaging as a Marker for Parkinson’s Disease. *Mov. Disord.* 2017, 32, 619–623.
- [5]. Deeraj Shetty, Kishor Rit, Sohail Shaikh, Nikita Patil, “Diabetes Disease Prediction Using Data Mining”, 2017 International Conference on Innovations in Information, Embedded and Communication Systems (ICIIECS)
- [6]. Samrat Kumar Dey, Ashraf Hossain, Md. Mahbubur Rahman, “Implementation of a Web Application to Predict Diabetes Disease: An Approach Using Machine Learning Algorithm”, 2018 21st International Conference of Computer and Information Technology (ICCIT)
- [7]. Dickson, D.W. Neuropathology of Parkinson disease. *Parkinsonism Relat. Disord.* 2018, 46 (Suppl. 1), S30–S33.

- [8]. Priyanka Sonar, Prof. K. JayaMalini, "Diabetes Prediction Using Different Machine Learning Approaches", Proceedings of the Third International Conference on Computing Methodologies and Communication (ICCMC 2019) IEEE Xplore Part Number: CFP19K25-ART; ISBN: 978-1-5386-7808-4
- [9]. International Diabetes Federation. IDF Diabetes Atlas, 9th edn. Brussels, Belgium: International Diabetes Federation, 2019.
- [10]. American Diabetes Association. Standards of medical care in diabetes-2020. Diabetes Care 2020; 43(Suppl. 1): S14–S31.
- [11]. Trends in coronary Heart Disease Epidemiology Center for Disease Control and Prevention (Heart Disease Facts).
- [12]. Asian Pacific Journal of Global Trend of Cancer Mortality rate: A 25-year study.
- [13]. International Diabetes Federation: Expenditure and deaths related to diabetes.
- [14]. Naveen Kishore G,V .Rajesh ,A.Vamsi Akki Reddy, K.Sumedh,T.rajesh Sai Reddy, "Prediction Of Diabetes Using Machine Learning Classification Algorithms".
- [15]. M.Marimuthu ,S.Deivarani ,R.Gayatri, "Analysis of Heart Disease Prediction using Machine Learning Techniques".
- [16]. Purushottam, Richa Sharma ,Dr. Kanak Saxena, "Efficient Heart Disease Prediction System".
- [17]. Adil Hussain She, Dr. Pawan Kumar Chaurasia," A Review on Heart Disease Prediction using Machine Learning Techniques".
- [18]. Times Of India: Cancer cases upswing 10% in 4 years to 13.9 lakh.
- [19]. Epidemiology of Diabetes :A report of Indian Heart Association.

AUTHORS PROFILE

1. Emad Naushad(1902220100066) Student, B.Tech Computer Science & Engineering, I.T.S Engineering College, Greater Noida.
2. Bhavishya Raj(1902220100056) Student, B.Tech Computer Science & Engineering, I.T.S Engineering College, Greater Noida.
3. Arpit Nirvan(1902220100038) Student, B.Tech Computer Science & Engineering, I.T.S Engineering College, Greater Noida.
4. Vrinda Sachdeva Associate Professor of Department of Computer Science & Engineering at I.T.S Engineering College, Greater Noida.