# D.Y. PATIL COLLEGE OF ENGINEERING &TECHNOLOGY, KASABA BAWADA, KOLHAPUR

## DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
## (ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING)



A Project-I Report

On

**"Heart Disease Prediction"**

Submitted by:

| Sr No. | Name | Roll no. |
|--------|------|----------|
| 1. | Bhagyashree M. Swami | 17 |
| 2. | Rohit V. Chitari | 09 |
| 3. | Prathmesh S. Suryavanshi | 15 |
| 4. | Aniket C. More | 14 |

Under the guidance of

Mr. Nitish M. Shinde

**Class :- TY BTech**

**(CSE-AIML)**

**Academic Year - (2023-2024)**

# D.Y. PATIL COLLEGE OF ENGINEERING &TECHNOLOGY, KASABA BAWADA, KOLHAPUR

## DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
## (ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING)

## CERTIFICATE

This is to certify that the following members have satisfactorily completed the Project-I work entitled **"Heart Disease Prediction"** at T.Y. B.Tech. CSE(AIML) Semester V prescribed in the curriculum of DYPCET Autonomy for theacademic year 2023-24.

| Roll no. | Name | Signature |
|---|---|---|
| 17 | Bhagyashree M. Swami | |
| 09 | Rohit V. Chitari | |
| 15 | Prathmesh S. Suryavanshi | |
| 14 | Aniket C. More | |

**Date**:

**Place**: Kolhapur

**Mr.  Nitish M. Shinde**                                                        **Dr. S. V. Patil**

(Project-I Guide)                                                                (HOD  CSE (AI-ML))

**Mr.  Nitish M. Shinde**                                                        **Dr. S. D.  Chede**

(Project-I Coordinator)                                                         (Principal)

**External Examiner**

# ACKNOWLEDGMENT

The success & find outcome of this project required a lot of guidance & assistance from many people and We are extremely privileged to have got all along the completion of our project. All that we have done is only due to such supervision & assistance & we would not forget to thanks them.

We owe our deep gratitude to our project guide and coordinator **Mr. Nitish M. Shinde** who took keen interest on our project work & guided us all along, till the completion of our project work by providing all the necessary Information for developing a good system.

We would like to express our gratitude and deep regards to department HOD **Dr. S. V. Patil** for his guidance and support throughout the completion of project work. We would like to express our heart full gratitude to him for his continuous encouragement & motivation.

It is our pleasure to acknowledge the help we have been received from institute and the Individual. We would like to thank our Principal **Dr. S. D. Chede** Principal, D. Y. Patil College of Engineering and Technology, Kolhapur in particular for always giving encouragement, support and the excellent facilities provided.

**Date:**

**Place:** Kolhapur

| Roll no. | Name | Signature |
|----------|------|-----------|
| 17 | Bhagyashree M. Swami | |
| 09 | Rohit V. Chitari | |
| 15 | Prathmesh S. Suryavanshi | |
| 14 | Aniket C. More | |

# ABSTRACT

Heart diseases remain a significant global health concern, contributing to a substantial number of morbidity and mortality cases each year. Early and accurate prediction of heart disease can aid in timely interventions and personalized treatments. This study explores the application of machine learning techniques for heart disease prediction using a comprehensive dataset of clinical and diagnostic attributes. The findings of this study indicate that machine learning algorithms can effectively predict heart disease based on the provided dataset. Certain algorithms demonstrate higher predictive accuracy and are particularly effective in identifying patients at risk. The study underscores the importance of feature selection in improving model performance and interpretability.

However, it is acknowledged that this research is limited by the nature and size of the dataset, and further validation using larger and diverse datasets is recommended. Additionally, the integration of expert medical knowledge and ethical considerations is crucial when deploying predictive models in real-world clinical settings. In conclusion, this study showcases the potential of machine learning in enhancing heart disease prediction accuracy. The results suggest that these techniques can complement traditional clinical assessment methods, enabling early intervention and personalized care for individuals at risk of Heart diseases.

# INDEX

# CHAPTER 1. INTRODUCTION

Heart disease is one of the leading causes of death worldwide, and early detectionplays a vital role in improving patient outcomes. With advancements in technology and the availability of large-scale medical data, predictive modeling has become a promising approach to identify individuals at risk of developing heart disease. In this context, a Heart Disease Prediction Program is designed to utilize machine learning algorithms and patient data to predict the likelihood of an individual developing heart disease in the future. The primary objective of the Heart Disease Prediction Program is to assist healthcare professionals in making informed decisions and providing timely interventions to high -risk patients. By analyzing various risk factors and patterns from patient data, the program aims to generate accurate predictions, thereby enabling early intervention, personalized treatment plans, and lifestyle modifications.

The domain information of a heart disease prediction program refers to the specific aspects related to the field or area in which the program operates. In the case of a heart disease prediction program, the domain information includes:

**Cardiovascular Medicine:** The program operates within the domain of cardiovascular medicine, focusing on the prediction and prevention of heart disease. It considers various cardiovascular conditions, including coronary artery disease, heart failure, arrhythmias, and valvular heart disease.

**Risk Factors:** The program takes into account a range of risk factors associated with heart disease. These factors may include age, gender, family history, smoking habits, high blood pressure, high cholesterol levels, diabetes, obesity, physical inactivity, and stress.

**Medical Data:** The program utilizes medical data to predict the likelihood of heart disease. This data includes patient information such as medical history, lifestyle factors, symptoms,and results from medical tests and diagnostic procedures. Examples of relevant tests include blood pressure measurements, cholesterol profiles, ECG recordings, stress tests, and echocardiograms.

**Machine Learning and Data Analysis:** The program may incorporate machine learning algorithms and data analysis techniques to analyze the input data and identify patterns or correlations that are indicative of heart disease. These algorithms can be trained on large datasets of patient information to improve the accuracy of predictions.

**Prevention and Treatment:** The program aims to provide insights into the risk of developing heart disease, allowing healthcare providers to implement appropriate preventive measures and treatments. This may include lifestyle modifications (e.g., diet, exercise), medication management, and referral to specialists for further evaluation or intervention.

# CHAPTER 2. LITERATURE REVIEW

**[1]** S. Ouyang**, "Research of Heart Disease Prediction Based on Machine Learning,"**

The use of massive clinical data in the medical field for supporting medical decision support is an inevitable development trend. Medical decision support is based on a variety of data sources accumulated and acquired in real-time in the clinic, and various machine learning algorithms are used to achieve classification of patient disease types or prediction of disease risks. This paper assists in performing cardiac disease prediction starting from different heart disease types (coronary heart disease) and data sets, summarizing the currently adopted machine learning diagnosis and prediction methods, highlighting the characteristics and differences of these methods, and analyzing the challenges and future developments. The results show that machine learning techniques have a wide range of applications in cardiac diseases. However, each machine learning method can only be applied to a specific scope due to the non-uniformity of medical data. At the end of the article, the prediction of heart disease is summarized.

**[2] Soni, Jyoti, Ujma Ansari, Dipesh Sharma, and Sunita Soni. "Predictive data mining for medical diagnosis: An overview of heart disease prediction."**

Medical data mining has great potential for exploring the hidden patterns in the data sets of the medical domain. These patterns can be utilized for clinical diagnosis. However, the available raw medical data are widely distributed, heterogeneous in nature, and voluminous. These data need to be collected in an organized form. This collected data can be then integrated to form a hospital information system. Data mining technology provides a useroriented approach to novel and hidden patterns in the data.

**[3] "HEART DISEASE PREDICTION SYSTEM"** Mrs.Jayashree L K, Sushmita Makapur, T D Vyshnavi, T D Prathyusha, V K Kavya

We referred about data analysis from this research paper along with the data flow diagrams.

# CHAPTER 3. PROBLEM STATEMENT

The high prevalence of heart disease and its significant impact on public health necessitate effective methods for early detection and risk assessment. Despite advancements in medical science, accurate and timely identification of individuals at risk of heart disease remains a challenge. This problem statement aims to develop a robust and reliable heart disease prediction model that leverages machine learning techniques to analyze comprehensive patient data and provide accurate risk assessment, facilitating proactive medical interventions and contributing to improved cardiovascular health outcomes.

The goal of this project is to develop a robust machine learning model for predicting the likelihood of heart disease in individuals based on a set of medical and lifestyle features. The model should analyze and interpret relevant data to provide accurate predictions and assist healthcare providers in making informed decisions about patient care and interventions. A comprehensive dataset containing a variety of attributes, including medical history, vital signs, blood tests, lifestyle choices, and possibly genetic information, will be provided. The dataset will be split into training and testing subsets to ensure unbiased model evaluation.

# CHAPTER 4. OBJECTIVES

**The objectives of the proposed work are as follows:**

**1.** Process the Dataset

**2.** Train the model using Logistic regression Algorithm

**3.** Test the dataset

**4.** Check the accuracy

**5.** Compare with existing system

**6.** Deployment

**7.** To provide effective methods of Heart Disease Prediction to decision makers towards better heart health management.

**8.** To provide recommendation to Heart health management system.

**9.** To prevent heart disease.

**10.** To make people aware about a healthy life.

**11.** To predict the seriousness of disease.

# CHAPTER 5. PROPOSED SYSTEM

## 5.1 System Requirements:

### 1. Hardware Requirements:

1. Processor: Intel Core i3
2. RAM: 8 GB
3. ROM: 256 GB

### 2. Software Requirements:

1. Operating System: Windows OS
2. Programming languages: Python
3. Database: Kaggle
4. Internet and Browsing Facilities
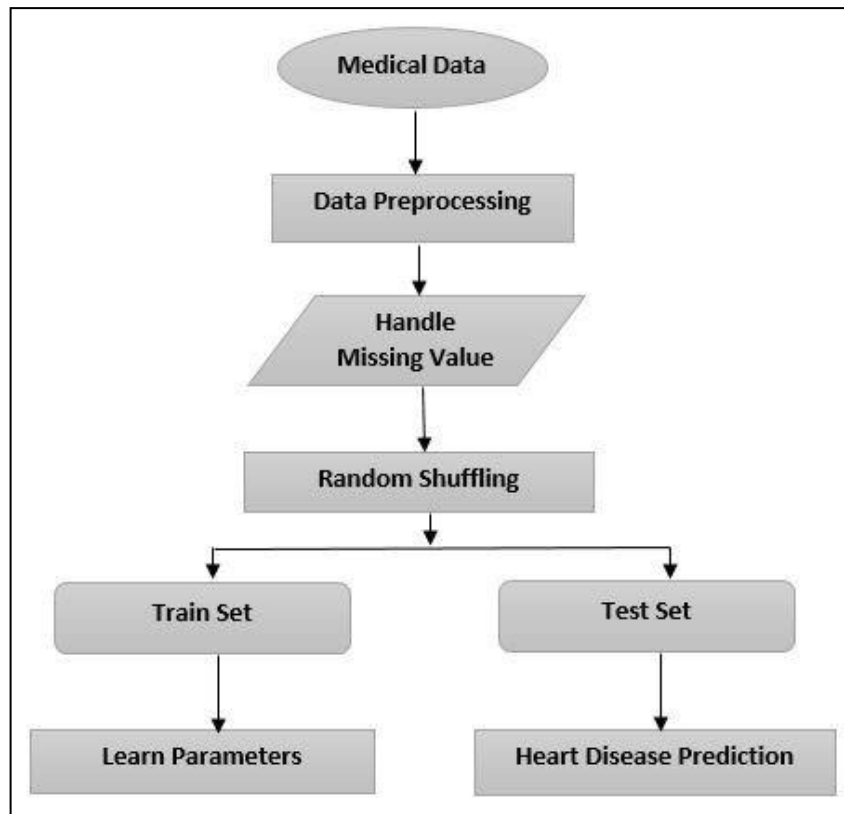
## 5.2       Architecture Diagram:



**Fig 5.2.1 Proposed System Architecture**

## 5.3    MODULES:

1. **Data Collection:** Gather a dataset containing relevant information about individuals, including their medical history, lifestyle factors, physiological measurements, and any diagnosed heart conditions. Ensure the dataset is diverse and representative of the population you aim to predict heart disease for.

2. **Data Preprocessing:** Clean the dataset by handling missing values, outliers, and inconsistencies. Perform feature engineering to extract meaningful features from the available data. This may involve transforming variables, normalizing data, or creating new features through domain knowledge.

3. **Data Split:** Divide the dataset into two or three subsets: a training set, a validation set, and optionally a test set. The training set is used to train the model, the validation set helps fine-tune the model's hyperparameters, and the test set provides an unbiased evaluation of the final model's performance.

4. **Model Selection:** Choose an appropriate machine learning algorithm for heart disease prediction. Some commonly used algorithms include logistic regression, decision trees, random forests, support vector machines (SVM), or neural networks. The choice of model depends on the dataset size, complexity, interpretability, and other specific requirements .

5. **Feature Selection:** Select the most informative features for heart disease prediction. This step can involve techniques like correlation analysis, recursive feature elimination, or regularization methods to identify the features that have the most significant impact onthe prediction task.

6. **Model Training:** Train the chosen model using the training dataset. The model learns from the input features and their corresponding heart disease labels.

7. **Model Evaluation:** Assess the performance of the trained model using the validation set . Common is Evalution metrics for classification tasks include accuracy, precision,recall F1 Score, and area under the Receiver operating characteristic curve (AUC-ROC).Adjust themodel's hyperparameters if necessary optimise  its performance.

8. **Model Testing:** Once satisfied with the model's performance, evaluate it on the independent test set. This provides an unbiased estimate of the model's generalization ability to predict heart disease on new, unseen data.

9. **Model Deployment:** If the model performs well on the test set, it can be deployed to make predictions on new data. The deployment can be as simple as providing a web interface or integrating the model into an existing healthcare system.

10. **Model Monitoring and Updating:** Continuously monitor the performance of the deployed model, and periodically update it to incorporate new data or account for changes in the population being predicted. This ensures that the model stays accurate and reliable over time.

# CHAPTER 6. EXPERIMENTAL WORK

1. **Data Collection and EDA:**
   - Data collection from website named Kaggle.
   - Performed cleaning on collected data, removed null values, removed outliers.
   - Feature selection for model training.

2. **Model Training:**
   - Divided data into dependent and independent features.
   - Train Test Split on data.
   - Trained model using different algorithms like Linear regression, Decision tree regression, Random Forest regression, Support Vector Regressor.

3. **Model Evolution:**
   - Continuously validate the model's performance.
   - Visualised various plots using seaborn and matplotlib.
   - Evaluated the model's performance using appropriate metrics, such as Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and correlation coefficients.

4. **Documentation:**
   - Maintain comprehensive documentation of data sources, models, algorithms, and system architecture.
   - Keep records of all changes and updates made during the experimental work.

# CHAPTER 7. RESULT ANALYSIS

## 1. Data Collection and EDA:

### 1.1 Data Collection:

- Read Dataset using pandas library.

```
In [2]: data = pd.read_csv("dataset heart.csv")
        data.head()
```

### 1.2 Remove Null Values:

**Checking Null Values**

```
In [8]: data.isnull().sum()

Out[8]: age        0
        sex        0
        cp         0
        trestbps   0
        chol       0
        fbs        0
        restecg    0
        thalach    0
        exang      0
        oldpeak    0
        slope      0
        ca         0
        thal       0
        target     0
        dtype: int64
```

```
In [9]: data.isnull().sum().sum()

Out[9]: 0
```

-From above observation we came to conclusion that there are no null values present in our dataset.

## 1.3 Description of complete Dataset:

```
4]:    data.describe()
```

| 4]: | | age | sex | cp | trestbps | chol | fbs | restecg | thalach | exang | oldpeak | slope | ca | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | count | 303.000000 | 303.000000 | 303.000000 | 303.000000 | 303.000000 | 303.000000 | 303.000000 | 303.000000 | 303.000000 | 303.000000 | 303.000000 | 303.000000 | 303.00 |
| | mean | 54.366337 | 0.683168 | 0.966997 | 131.623762 | 246.264026 | 0.148515 | 0.528053 | 149.646865 | 0.326733 | 1.039604 | 1.399340 | 0.729373 | 2.31 |
| | std | 9.082101 | 0.466011 | 1.032052 | 17.538143 | 51.830751 | 0.356198 | 0.525860 | 22.905161 | 0.469794 | 1.161075 | 0.616226 | 1.022606 | 0.61 |
| | min | 29.000000 | 0.000000 | 0.000000 | 94.000000 | 126.000000 | 0.000000 | 0.000000 | 71.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.00 |
| | 25% | 47.500000 | 0.000000 | 0.000000 | 120.000000 | 211.000000 | 0.000000 | 0.000000 | 133.500000 | 0.000000 | 0.000000 | 1.000000 | 0.000000 | 2.00 |
| | 50% | 55.000000 | 1.000000 | 1.000000 | 130.000000 | 240.000000 | 0.000000 | 1.000000 | 153.000000 | 0.000000 | 0.800000 | 1.000000 | 0.000000 | 2.00 |
| | 75% | 61.000000 | 1.000000 | 2.000000 | 140.000000 | 274.500000 | 0.000000 | 1.000000 | 166.000000 | 1.000000 | 1.600000 | 2.000000 | 1.000000 | 3.00 |
| | max | 77.000000 | 1.000000 | 3.000000 | 200.000000 | 564.000000 | 1.000000 | 2.000000 | 202.000000 | 1.000000 | 6.200000 | 2.000000 | 4.000000 | 3.00 |

## 1.4 Feature selection for model training:

```
f, ax = plt.subplots(figsize=(15, 10))
sns.set(font_scale=1.5)
hm = sns.heatmap(cm,
                cbar=True,
                annot=True,
                square=True,
                fmt='.2f',
                annot_kws={'size': 15},
                yticklabels=cols,
                xticklabels=cols)
plt.show()
```

## 2. Model Training:

### 2.1 Divide data into independent and dependent features:

```
In [15]: X=data[['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs', 'restecg', 'thalach',
            'exang', 'oldpeak', 'slope', 'ca', 'thal', ]].values
         y=data[['target']].values
```

- Variable x denotes the independent features.
- Variable y denotes dependent feature.

### 2.2 Train Test Split:

```
In [16]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)
```

- For splitting the data into train set and test set we imported train_test_split from sklearn model selection.
- Test size of dataset is 0.30 and training size is 0.70.
- When we use random state fixed examples are taken for training and testing.

### 2.3 K Nearest Neighbors Regressor:

**KNN**

```
In [23]: knn_model = KNeighborsClassifier(n_neighbors=5)

         knn_model.fit(X_train, y_train.ravel())

Out[23]: KNeighborsClassifier()

In [24]: knn_train_accuracy = knn_model.score(X_train, y_train)
         print('KNN Training Accuracy: %.2f' % knn_train_accuracy)

         KNN Training Accuracy: 0.78

In [25]: knn_test_accuracy = knn_model.score(X_test, y_test)
         print('KNN Test Accuracy: %.2f' % knn_test_accuracy)

         KNN Test Accuracy: 0.67
```

## 2.4 Support Vector Regressor:

### SVM

```
In [29]: svm_model = SVC(kernel='linear', C=1.0, random_state=0)

         svm_model.fit(X_train, y_train.ravel())

Out[29]: SVC(kernel='linear', random_state=0)
```

```
In [30]: svm_train_accuracy = svm_model.score(X_train, y_train)
         print('SVM Training Accuracy: %.2f' % svm_train_accuracy)

         SVM Training Accuracy: 0.85
```

```
In [31]: svm_test_accuracy = svm_model.score(X_test, y_test)
         print('SVM Test Accuracy: %.2f' % svm_test_accuracy)

         SVM Test Accuracy: 0.81
```

## 2.5 Decision Tree Regressor:

### Decision Tree

```
In [26]: dt_model = DecisionTreeClassifier(random_state=0)

         dt_model.fit(X_train, y_train.ravel())

Out[26]: DecisionTreeClassifier(random_state=0)
```

```
In [27]: dt_train_accuracy = dt_model.score(X_train, y_train)
         print('Decision Tree Training Accuracy: %.2f' % dt_train_accuracy)

         Decision Tree Training Accuracy: 1.00
```

```
In [28]: dt_test_accuracy = dt_model.score(X_test, y_test)
         print('Decision Tree Test Accuracy: %.2f' % dt_test_accuracy)

         Decision Tree Test Accuracy: 0.75
```

## 2.6 Random Forest Regressor:

### Random forest

```
In [20]: rf_model = RandomForestClassifier(n_estimators=100, random_state=0)

rf_model.fit(X_train, y_train.ravel())

Out[20]: RandomForestClassifier(random_state=0)

In [21]: rf_train_accuracy = rf_model.score(X_train, y_train)
print('Random Forest Training Accuracy: %.2f' % rf_train_accuracy)

Random Forest Training Accuracy: 1.00

In [22]: rf_test_accuracy = rf_model.score(X_test, y_test)
print('Random Forest Test Accuracy: %.2f' % rf_test_accuracy)

Random Forest Test Accuracy: 0.84
```

## 2.7 Logistic Regression:

### Logistic Regression

```
In [16]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)

In [17]: model = LogisticRegression( C=200, penalty='l1',solver='liblinear')

model.fit(X_train, y_train)

C:\Users\hp\anaconda3\lib\site-packages\sklearn\utils\validation.py:63: DataConversionWarning: A
n a 1d array was expected. Please change the shape of y to (n_samples, ), for example using rave
  return f(*args, **kwargs)

Out[17]: LogisticRegression(C=200, penalty='l1', solver='liblinear')

In [18]: model.score(X_train, y_train)

Out[18]: 0.8679245283018868

In [19]: print ('Training Accuracy: %.2f'%model.score(X_train,y_train))
print ('Test Accuracy: %.2f' % model.score(X_test,y_test))

Training Accuracy: 0.87
Test Accuracy: 0.81
```

## 3. Model Evalution:

### 3.1 Performance Table:

| Sr. No. | Regression Model Name | Accuracy |
|---------|------------------------|----------|
| 1. | KNNeighbors Regressor | 78.00% |
| 2. | Support Vector Regression | 85.00% |
| 3. | Decision Tree Regression | 100.00% |
| 4. | Random Forest Regression | 100% |
| 5. | Logistic Regression | 87.00% |

- In above table we have accuracy of all the regression models.
- Besides Decision Tree Algorithms gives the maximum accuracy.
- Model also worked well on new data and predicted nearby correct output.

**There are a few reasons why decision trees and random forests are well-suited for heart disease prediction:**

1. **They are interpretable**: Unlike some other machine learning algorithms, decision trees and random forests are relatively easy to interpret. This means that it is possible to understand how the algorithm is making its predictions, which can be helpful for identifying important risk factors for heart disease.

2. **They are robust to noise and outliers**: Decision trees and random forests are able to handle noisy data and outliers relatively well. This is important for heart disease prediction, as the data can be noisy due to factors such as measurement errors and incomplete medical records.

3. **They are efficient to train and predict with**: Decision trees and random forests can be trained and used to make predictions efficiently, even on large datasets. This is important for heart disease prediction, as it allows the algorithm to be used in real-time applications, such as screening patients for heart disease risk.

### 3.2 Prediction on new data:

## Logistic Regression

```
In [16]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)

In [17]: model = LogisticRegression( C=200, penalty='l1',solver='liblinear')

         model.fit(X_train, y_train)

         C:\Users\hp\anaconda3\lib\site-packages\sklearn\utils\validation.py:63: DataConversionWarning: A
         n a 1d array was expected. Please change the shape of y to (n_samples, ), for example using rave
           return f(*args, **kwargs)

Out[17]: LogisticRegression(C=200, penalty='l1', solver='liblinear')

In [18]: model.score(X_train, y_train)

Out[18]: 0.8679245283018868

In [19]: print ('Training Accuracy: %.2f'%model.score(X_train,y_train))
         print ('Test Accuracy: %.2f' % model.score(X_test,y_test))

         Training Accuracy: 0.87
         Test Accuracy: 0.81
```

# CHAPTER 8: CONCLUSION

In conclusion, the development and implementation of machine learning algorithms for heart disease prediction have shown promising results. These models leverage various data sources, such as patient demographics, medical history, and diagnostic tests, to make accurate predictions about an individual's risk of heart disease. The use of techniques like logistic regression, decision trees, random forests, support vector machines, and deep learning has demonstrated their ability to effectively classify and predict heart disease. These models have the potential to assist healthcare professionals in identifying at-risk individuals, enabling timely interventions and personalized treatment plans. However, further research and validation are needed to ensure the models' accuracy, generalizability, and real-world applicability, while addressing potential ethical and privacy concerns associated with the use of sensitive medical data. Overall, the development of heart disease prediction models using machine learning is a promising avenue in improving early detection and prevention of Heart Disease Prediction

# REFERENCES

[1] S. Ouyang, "Research of Heart Disease Prediction Based on Machine Learning," 2022 5th International Conference on Advanced Electronic Materials, Computers and Software Engineering (AEMCSE), Wuhan, China, 2022, pp. 315-319, doi:10.1109/AEMCSE55572.2022.00071.
https://ieeexplore.ieee.org/abstract/document/9948280/

[2] Soni, Jyoti, Ujma Ansari, Dipesh Sharma, and Sunita Soni. "Predictive data mining for medical diagnosis: An overview of heart disease prediction." International Journal of Computer Applications 17, no. 8 (2011): 43-48.
https://www.academia.edu/download/79534142/5a18f6653b56138cd5196d20e2f39de189e3.pdf

[3] "HEART DISEASE PREDICTION SYSTEM" Mrs.Jayashree L K, Sushmita Makapur, T D Vyshnavi, T D Prathyusha, V K Kavya http://www.irjcs.com/

External Links:

https://www.researchgate.net/
https://towardsdatascience.com/
https://scholar.google.com/

**Date:**

**Place:** Kolhapur

Mr. Nitish Shinde                                          Dr. S. V. Patil

(Project- I Guide & Coordinator)                    (H.O.D , CSE AI-ML)