# INTERPRETABLE MACHINE LEARNING MODELS FOR HEART DISEASE DETECTION

[1]Mr.S. Sathish Kumar, [2]B. Mamatha, [3]E. Swarna, [4]E. Srinivas, [5]N. Jayanth Reddy

[1]Assistant Professor, [2,3,4,5]Student
[1]Department Of Computer Science and Engineering,
[1]JB INSTITUTE OF ENGINEERING AND TECHNOLOGY
(Affiliated to Jawaharlal Nehru Technological University, Hyderabad, India)

*Abstract:* Heart disease is a type of disease that affects the heart or blood vessels. The risk of certain heart disease may increase with smoking, high blood pressure, high cholesterol, poor diet, lack of exercise, and obesity. The most common heart disease is coronary artery disease that can lead to chest pain, heart disease, or stroke. Early detection of heart disease is important for saving lives. In this paper, we explored various ways to learn the machine in predicting heart disease. But the biggest problem in the working environment is the adoption of black box machine learning models. Since doctors often give diagnoses based on their experience and knowledge-based thinking, it becomes difficult for them to accept vague and difficult-to-understand models in diagnosing a serious illness involving the cost of human life. So, to deal with their dubiety we have designed an interpretable machine learning model which presents key attributes/features and their range of values responsible in predicting class type of heart disease.

*Keywords* - **Random Forest, AdaBoost, XGBoost, Light GBM, Multilayer Perceptron, Lime, eli5.**

## 1. INTRODUCTION

Heart disease is very common these days, describing many conditions that can affect your heart. The World Health Organization estimates that 17.9 million people worldwide die from heart disease. These conditions are aggravated by smoking, excessive drinking, lack of exercise, and insomnia. Many health conditions, your lifestyle, your age and family history can increase your risk of heart disease and stroke. A heart attack is the leading cause of death worldwide. Lifestyle changes will be an important factor in reducing risk and it is expected that the development of calculations that can reduce heart disease will significantly reduce cardiovascular mortality and early detection could lead to significant reductions in health care costs. An estimated 17.9 million people died of CVDs in 2019, which means 32% of all deaths in the world. Of these deaths, 85% were caused by heart disease and stroke. The use of analysis in health care enhances care by performing preventive care and helps us to have a clear study and awareness of the causes of heart disease. It is the leading cause of death among the elderly.

Our paper can help predict people who may be diagnosed with heart disease with the help of their medical history. It recognizes who all have any symptoms of heart disease such as chest pain or high blood pressure and can help diagnose the disease with a little medical examination and effective treatment, so that they can be treated appropriately. More information and symptoms related to heart disease i.e., age, blood pressure, cholesterol, hyper tension etc.

The set of cardiovascular data basically contains the above-mentioned information as well as summarized and collected data from patients. With the growing number of patients and data sets we need to use an effective approach that can assist the physician in providing accurate predictors of the patient's heart disease.

The purpose of this paper is to assess whether a patient is likely to be diagnosed with cardiovascular disease based on their medical characteristics such as gender, age, chest pain, fasting blood sugar level, etc. The database is selected from the UCI database with the patient's medical history and qualifications. Using this database, we predict whether a patient may have heart disease or not. To predict this, we use 14 patient medical features and differentiate if the patient is likely to have heart disease. These medical attributes are trained under five algorithms. We classify patients who are at risk for heart disease or not and this method is completely economical.

## 2. RELATED WORK

K.Prasanna Lakshmi, Dr. C.R.K.Reddy (2015) designed "Fast Rule-BasedHeart Disease Prediction using Associative Classification Mining". In the proposed Stream Associative Classification Heart Disease Prediction (SACHDP), we used associative classification mining over landmark window of data streams. This paper contains two phases: one is generating rules from associative classification mining and next one is pruning the rules using chi- square testing and arranging the rules in an order to form a classifier. Using these phase to predict the heart disease easily [1].

M.Satish, et al. (2015) used different Data Mining techniques like Rule based, Decision Tree, Navie Bayes, and Artificial Neural Network. An efficient approach called pruning classification association rule (PCAR) was used to generate association rules from cardiovascular disease warehouse for predictionof Heart Disease. Heart attack data warehouse was used for pre-processing for mining. All the above discussed data mining technique were described [2].

Lokanath Sarangi, Mihir Narayan Mohanty, Srikanta Pattnaik (2015) "AnIntelligent Decision Support System for Cardiac Disease Detection", designed a cost efficient model by using genetic algorithm optimizer technique. The weights were optimized and fed as an input to the given network. The accuracyachieved was 90% by using the hybrid technique of GA and neural networks [3].

Ashir Javeed, Shijie Zhou et al. (2017) designed "An Intelligent Learning System based on Random Search Algorithm and Optimized Random Forest Model for Improved Heart Disease Detection". This paper uses random search algorithm (RSA) for factor selection and random forest model for diagnosing the cardiovascular disease. This model is principally optimized for using grid search algorithmic program. Two forms of experiments are used forcardiovascular disease prediction. In the first form, only random forest model is developed and within the second experiment the proposed Random Search Algorithm based random forest model is developed. This methodology isefficient and less complex than conventional random forest model. Comparing to conventional random forest it produces 3.3% higher accuracy. The proposedlearning system can help the physicians to improve the quality of heart failure detection [4].

Bo Jin, Chao Che et al. (2018) proposed a "Predicting the Risk of Heart FailureWith EHR Sequential Data Modeling" model designed by applying neural network. This paper used the electronic health record (EHR) data from real- world datasets related to congestive heart disease to perform the experiment andpredict the heart disease before itself. We tend to used one-hot encryption and word vectors to model the diagnosing events and foretold coronary failure events victimization the essential principles of an extended memory network model. By analyzing the results, we tend to reveal the importance of respectingthe sequential nature of clinical records [5].

Aakash Chauhan et al. (2018) presented "Heart Disease Prediction using Evolutionary Rule Learning". This study eliminates the manual task that additionally helps in extracting the information (data) directly from the electronic records. To generate strong association rules, we have applied frequent pattern growth association mining on patient's dataset. This will facilitate (help) in decreasing the amount of services and shown that overwhelming majority of the rules helps within the best prediction of coronarysickness [6].

"Prediction and Diagnosis of Heart Disease by Data Mining Techniques"designed by Boshra Bahrami, Mirsaeid Hosseini Shirvani. This paper uses various classification methodology for diagnosing cardiovascular disease. Classifiers like KNN, SVO classifier and Decision Tree are used to divide the datasets. Once the classification and performance evaluation the Decision tree is examined as the best one for cardiovascular disease prediction from the dataset [7].

Mamatha Alex P and Shaicy P Shaji (2019) designed "Prediction and Diagnosisof Heart Disease Patients using Data Mining Technique". This paper uses techniques of Artificial Neural Network, KNN, Random Forest and Support Vector Machine. Comparing with the above-mentioned classification techniquesin data mining to predict the higher accuracy for diagnosing the heart disease isArtificial Neural Network [8].

## 3. MATERIALS AND METHODS

**3.1 DATA COLLECTION:** In the UCI archive a set of cardiovascular data is extracted. It basically contains the most commonly used traits in diagnosing heart disease in individuals. The main goal is to detect the presence of heart disease. The total data contains all medical results from various forums added to form a data set, consisting of 14 columns and 1026 rows.

**3.2 TESTING AND TRAINING:** In Machine Education, we usually divide our data into two sub-sets: training data and test data we call Test / Train Division. In this paper we divide our training data into 80 and test as 20.

## 3.3 ALGORITHMS:

**3.3.1 Random Forest (RF):** Random Forest planning is well known as a combination of the combination used in the field of mechanical and data science in various areas of the system. This method uses a "coherent combination" that suits the same number of decision dividers, in small samples of the data set and uses multiple or intermediate voting results or end results. It therefore reduces the problem of over-equality and increases the accuracy and control of forecasting. Therefore, the RF learning model with multiple decision trees is generally more accurate than a single tree-based model. To create a series of decision trees with a controlled variety, it includes a combination of bootstrap (bags) and a selection of random feature. It adapts to both planning and retrospective problems and fits well in both class and continuous values.

**3.1.2 Adaptive Boosting (AdaBoost):** Adaptive Boosting (AdaBoost) is an integrated learning process that uses a recurring approach to empower poor designers by learning from their mistakes. This was developed by Yoav Freund et al. also known as "meta-learning". Unlike a random forest that uses the same combination, Adaboost uses a "sequential integration". It creates a powerful classification by combining multiple subdivisional dividers to obtain a good high precision separator. In that sense, AdaBoost is called an adaptive classifier by greatly improving the efficiency of the separator, but in some cases, it can cause overcrowding. AdaBoost is best used to maximize the performance of decision trees, the basic rate, in binary separation problems, however, it is sensitive to sound and external data.

**3.1.3 Extreme gradient boosting (XGBoost):** Gradient Boosting, like Random Forests above, is an integrated learning algorithm that produces a final model based on a series of individual models, usually cutting trees. A gradient is used to reduce weight loss, such as how neural networks use gradient reduction to prepare weights. Extreme Gradient Boosting (XGBoost) is a type of gradient extension that takes more detailed measurements to be considered when determining the best model. It calculates the gradients of the second order of loss work to reduce losses and improved familiarity (L1 and L2), which reduces excessive balance, and improves model performance and performance. XGBoost is fast translation and can handle large data sets.

**3.1.4 Light Gradient Boosting (LightGBM):** LightGBM is a variation in the gradient expansion proposed by Ke et al. in 2017. Gradient boosting refers to an ensemble model based on the decision tree as a weak student. The guessing power and calculation costs of this algorithm deteriorate if a large amount of data is obtained, or the size of the attribute is high. The LightGBM model can overcome these limitations by using a one-sided sample based on gradient (GOSS) and special feature integration (EFB) techniques. Moreover, the LightGBM model grows its trees using a clever leaf strategy, rather than a clever tree method. This strategy elevates trees upside down, while other algorithms grow horizontally.

**3.1.5 Multilayer Perceptron (MLP):** MLP models are subdivisions of ANN relay that contain input layer, output layer, and one or more hidden layers. The model starts by transmitting the signal forward from the input layer to the hidden layer and finally to the output layer. Continuously, the error signal is streamed back to the input layer. The learning algorithm adjusts network weights and biases until the error reaches an acceptable level. In the MLP model with a single hidden layer, advanced parameters include (1) (activation): hidden layer activation function, (2) (solution): learning algorithm, and (3) (hidden_layers_sizes): number of neurons in the layer hidden.

## 3.4 INTERPRETABLE MODELS FOR MACHINE LEARNING ALGORITHMS

**Permutation Importance (Eli5):** It a provides a way to compute feature importance's for any black box estimator by measuring how score decreases when a feature is not available.

**Local Interpretable Model-agnostic Explanations (LIME):** LIME takes the interpretive representation of these sample points, determines their assumptions and builds a linear weight by reducing losses and complexity. Spatial definition method LIME interprets each prediction by reading the local translation model. The feeling behind LIME is that a sample of events in near and far areas is the interpretive manifestation of the original input. Then LIME takes the interpretive representation of these sample points, determines their assumptions and builds a linear weight by reducing losses and complexity. The weight of the samples is based on their distance from the original position. Points weigh less as the points go farther. The description is reliable locally, which means it represents a predictive model for local events.
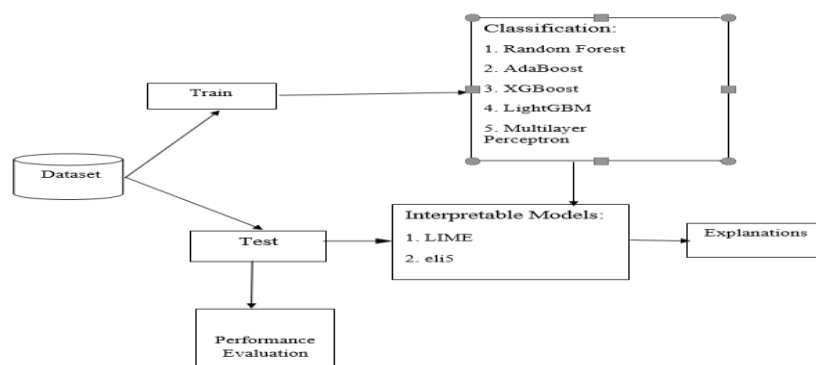
## 4. PROPOSED SYSTEM



Table 4.1 Proposed System

## 5. EXPERIMENTAL RESULTS AND DISCUSSIONS

This paper focuses on learning algorithms for five data machines namely: Random Forest Classifier, Adaptive Boosting (Ada boost), Extreme Gradient Boost (XG Boost), Gradient Light Development (Light GBM), -Multilayer Perceptron. We have also used two Translated Models in these algorithms, the Definitions of the Most Interpreted Model- agnostic (LIME) and Significance of Consent (Eli5) to make it clearer.

| Classifiers | Accuracy | F1 Score | Precision | Recall |
|---|---|---|---|---|
| Random Forest | 98.537 | 0.98605 | 0.97248 | 1.0 |
| XGBoost | 92.683 | 0.93088 | 0.90991 | 0.95283 |
| Ada Boost | 97.561 | 0.97696 | 0.95495 | 1.0 |
| Light GBM | 90.244 | 0.90909 | 0.87719 | 0.9434 |
| Multilayer perceptron | 89.756 | 0.90411 | 0.87611 | 0.93396 |

Table-5.1: Performance Evaluation of Classifiers

### 5.1 Local Interpretable Model-agnostic  Explanations
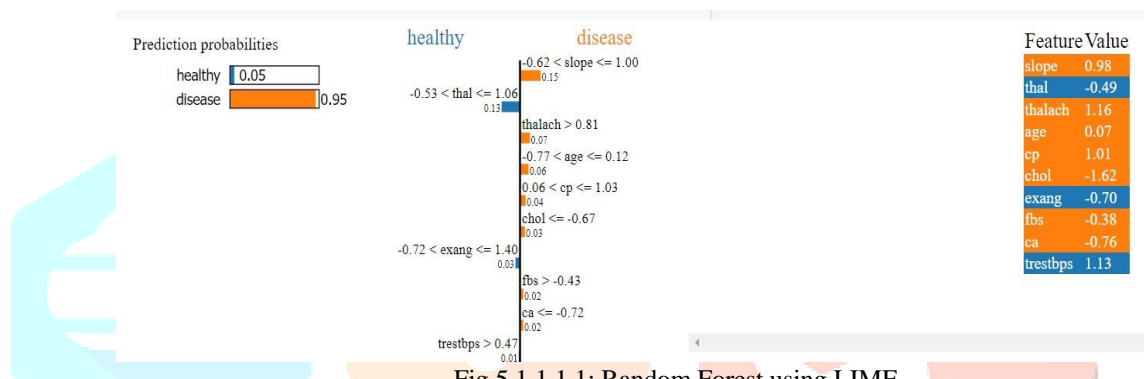
### 5.1.1 LIME:

### 5.1.1.1 Random Forest using LIME



Fig 5.1.1.1.1: Random Forest using LIME

We have used LIME in the Random Forest Classifier. The figure above shows that the Prediction Probalities of Healthy is 0.05 and Disease is 0.95. Common (healthy) attributes are thal, exang, trestbps and Qualities that lead to heart disease are Slope, thalach, age, cp, chol, fbs, ca.

### 5.1.1.2 AdaBoost using LIME:



Fig 5.1.1.1.2: AdaBoost using LIME

We have used LIME in the AdaBoost Algorithm. The figure above shows that the Prediction Probalities of Healthy is 0.03 and Disease is 0.97. Common (healthy) qualities are thal, exang and Qualities that lead to heart disease are Slope, thalach, age, cp, chol, fbs, ca, oldpeak.
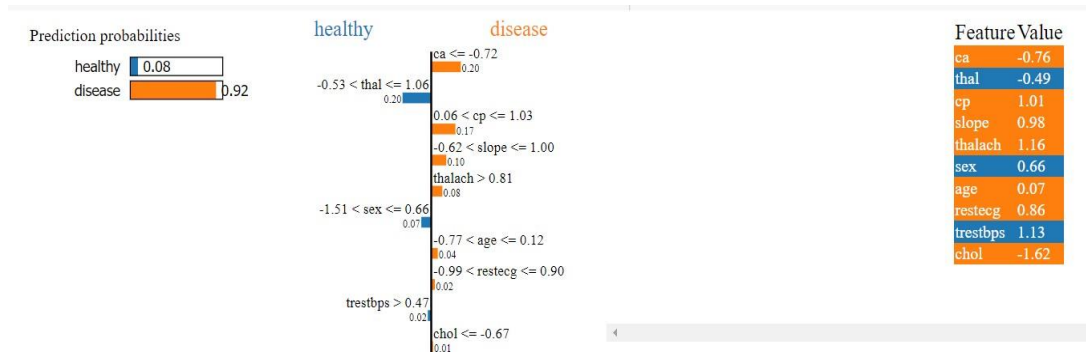
**5.1.1.3 LightGBM using LIME:**



Fig 5.1.1.1.3: LightGBM using LIME

We have applied LIME to the LightGBM Algorithm. The figure above shows that the probability of a healthy prediction is 0.08 and the Disease is 0.92. Common (healthy) attributes are thal, sex, trestbps and Qualities that lead to heart disease are ca, cp, slope, thalach, age, restecg, chol.
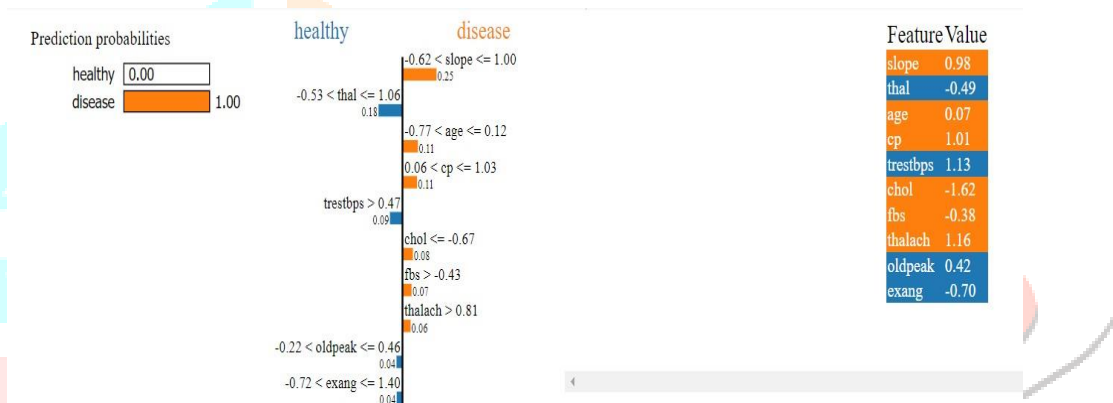
**5.1.1.4 XGBoost using LIME:**



Fig 5.1.1.1.4: XGBoost using LIME

We have used LIME in the XGBoost Algorithm. The figure above shows that Prediction Probalities of Healthy is 0.00 and Disease is 1.00. The right qualities normal (healthy) are thal, trestbps, oldpeak, exang and Qualities that lead to heart disease slope, age, cp, chol, fbs, thalach.
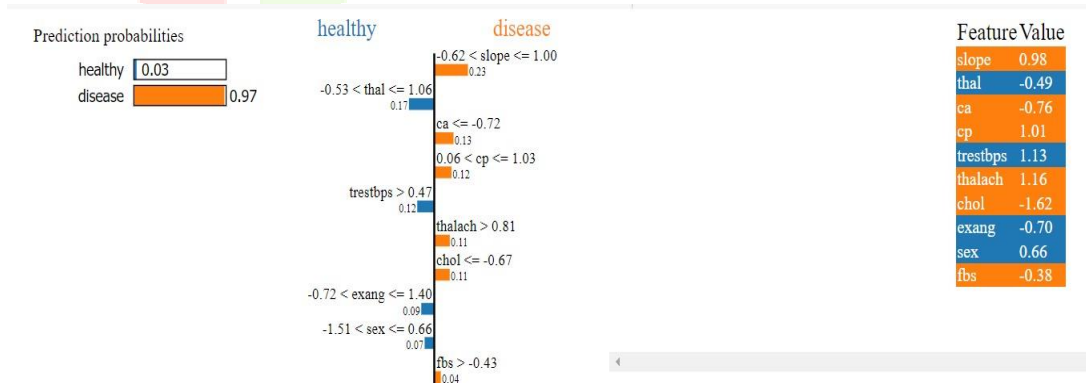
**5.1.1.5 Multilayer Perceptron using LIME:**



Fig 5.1.1.1.5: Multilayer Perceptron using LIME

We have used LIME in the Multlayer Perceptron Algorithm.The figure above shows that the Prediction Probalities of Healthy is 0.03 and the Disease is 0.97. Common (healthy) qualities are thal, trestbps, exang, sex and Qualities that lead to heart disease slope, ca, cp, thalach, chol, fbs.

### 5.1.2 Performance Importance(eli5):

Highly black figures are the most important factors, and those that look down with light shades are less important. The first number in each row indicates how much the performance of the model decreased by random shuffling (in this case, using "precision" as the performance metric). Random measurement in calculating the value of permits is done by repeating the process with multiple shifts. The number after ± measures how performance varies from one setup to the next.

### 5.1.2.1 Random Forest using eli5:

| Weight | Feature |
|---|---|
| 0.0985 ± 0.0285 | ca |
| 0.0722 ± 0.0143 | cp |
| 0.0585 ± 0.0138 | thal |
| 0.0371 ± 0.0132 | oldpeak |
| 0.0341 ± 0.0239 | thalach |
| 0.0273 ± 0.0228 | exang |
| 0.0273 ± 0.0048 | chol |
| 0.0263 ± 0.0236 | sex |
| 0.0205 ± 0.0114 | trestbps |
| 0.0127 ± 0.0099 | slope |
| 0.0107 ± 0.0129 | restecg |
| 0.0107 ± 0.0189 | age |
| 0.0020 ± 0.0048 | fbs |

Fig 5.1.2.1.1: Random Forest using eli5

The above Fig shows that ca, cp, thal, oldpeak and thalach are top 5 important features. The weight of "ca" is 0.0985 ± 0.0285. The model performance is decreased with a random shuffling is 0.0985 and The performance varied from one-reshuffling to the next is 0.0285.

**y=1 (probability 0.974) top features**

| Contribution? | Feature | Value |
|---|---|---|
| +0.514 | <BIAS> | 1.000 |
| +0.096 | thal | -0.486 |
| +0.094 | cp | 1.009 |
| +0.084 | ca | -0.763 |
| +0.048 | oldpeak | -0.962 |
| +0.047 | age | -1.969 |
| +0.044 | sex | -1.519 |
| +0.041 | exang | -0.697 |
| +0.021 | chol | -0.510 |
| +0.016 | thalach | 0.179 |
| +0.001 | fbs | -0.381 |
| -0.000 | restecg | 0.860 |
| -0.002 | trestbps | 0.419 |
| -0.029 | slope | -0.652 |

Fig 5.1.2.1.2 : Random Forest using eli5

To make random forest predictions more interpretable, every prediction of the model can be presented as a sum of feature contributions (plus the bias), showinghow the features lead to a particular prediction. In above plot, ELI5 does it by showing weights for each feature with their actual value depicting how influential it might have been in contributing to the final predictiondecisi across all trees. In the above individual prediction, the top 3 influential featuresseems to be, after the bias, cp and ca.

### 5.1.2.2 XGBoost using eli5:

| Weight | Feature |
|---|---|
| 0.1044 ± 0.0280 | ca |
| 0.0780 ± 0.0338 | cp |
| 0.0537 ± 0.0185 | sex |
| 0.0380 ± 0.0292 | oldpeak |
| 0.0351 ± 0.0156 | thal |
| 0.0312 ± 0.0181 | chol |
| 0.0293 ± 0.0276 | thalach |
| 0.0195 ± 0.0400 | age |
| 0.0185 ± 0.0114 | restecg |
| 0.0156 ± 0.0039 | slope |
| 0.0146 ± 0.0151 | trestbps |
| 0.0088 ± 0.0073 | fbs |
| -0.0049 ± 0.0138 | exang |

Fig 5.1.2.2.1: XG Boost using eli5

Fig. Above shows that ca, cp, sex, oldpeak and thal are 5 key factors. The weight of "ca" is 0.1044 ± 0.0280. Model performance is reduced by random shuffling by 0.1044 and Performance is split from one to next reversal by 0.0280.

**y=1** (probability **0.998**, score **6.407**) top features

| Contribution? | Feature | Value |
|---:|:---|---:|
| +1.399 | cp | 1.009 |
| +1.371 | sex | -1.519 |
| +1.365 | ca | -0.763 |
| +1.180 | age | -1.969 |
| +0.968 | thal | -0.486 |
| +0.570 | chol | -0.510 |
| +0.370 | oldpeak | -0.962 |
| +0.367 | exang | -0.697 |
| +0.144 | trestbps | 0.419 |
| +0.024 | <BIAS> | 1.000 |
| -0.097 | fbs | -0.381 |
| -0.191 | thalach | 0.179 |
| -0.348 | restecg | 0.860 |
| -0.714 | slope | -0.652 |

Fig 5.1.2.2.2: XG Boost using eli5

In the above plot, ELI5 does this by showing the weights of each element by their actual value which shows how much it would contribute to contributing to the final predictive decision in all trees. In each of the above predictions, the top 3 influential factors appear to be cp, gender and ca.
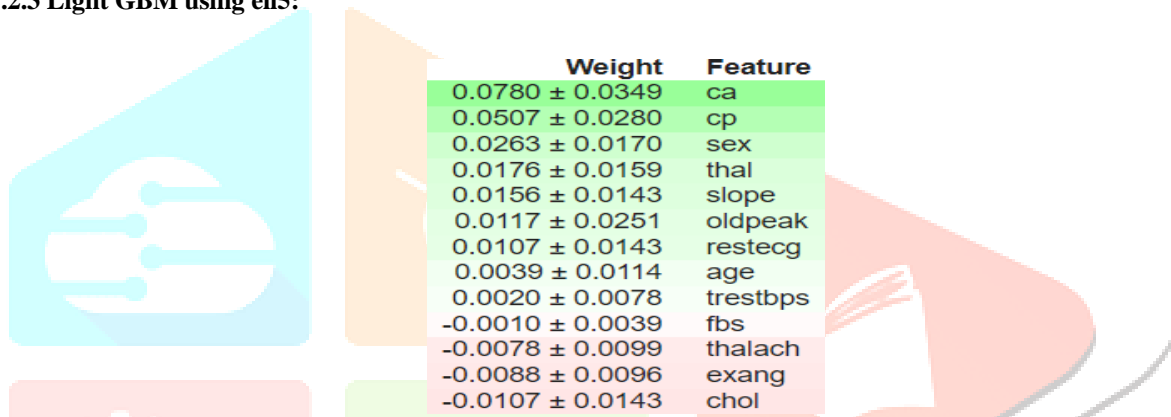
### 5.1.2.3 Light GBM using eli5:

| Weight | Feature |
|:---:|:---|
| 0.0780 ± 0.0349 | ca |
| 0.0507 ± 0.0280 | cp |
| 0.0263 ± 0.0170 | sex |
| 0.0176 ± 0.0159 | thal |
| 0.0156 ± 0.0143 | slope |
| 0.0117 ± 0.0251 | oldpeak |
| 0.0107 ± 0.0143 | restecg |
| 0.0039 ± 0.0114 | age |
| 0.0020 ± 0.0078 | trestbps |
| -0.0010 ± 0.0039 | fbs |
| -0.0078 ± 0.0099 | thalach |
| -0.0088 ± 0.0096 | exang |
| -0.0107 ± 0.0143 | chol |

Fig 5.1.2.3.1 LightGBM using eli5

The above Fig shows that ca, cp, sex, thal and slope are top 5 importantfeatures. The weight of "ca" is $0.0780 \pm 0.0349$. The model performance is decreased with a random shuffling is 0.0780 and The performance varied from one-reshuffling to the next is 0.0349.

| Contribution? | Feature | Value |
|---:|:---|---:|
| +0.644 | cp | 1.009 |
| +0.509 | ca | -0.763 |
| +0.418 | thal | -0.486 |
| +0.304 | age | -1.969 |
| +0.291 | sex | -1.519 |
| +0.207 | chol | -0.510 |
| +0.166 | exang | -0.697 |
| +0.161 | oldpeak | -0.962 |
| +0.064 | <BIAS> | 1.000 |
| +0.018 | restecg | 0.860 |
| -0.006 | thalach | 0.179 |
| -0.080 | trestbps | 0.419 |
| -0.206 | slope | -0.652 |

Fig 5.1.2.3.2 LightGBM using eli5

In above plot, ELI5 does it by showing weights for each feature with their actualvalue depicting how influential it might have been in contributing to the final prediction decision across all trees. In the above individual prediction, the top 3influential features seem to be, after the cp, ca and thal.

**5.1.2.4AdaBoost using eli5:**

| Weight | Feature |
|---|---|
| 0.0966 ± 0.0357 | ca |
| 0.0615 ± 0.0210 | thal |
| 0.0615 ± 0.0170 | cp |
| 0.0390 ± 0.0214 | oldpeak |
| 0.0234 ± 0.0272 | thalach |
| 0.0215 ± 0.0236 | sex |
| 0.0195 ± 0.0247 | exang |
| 0.0185 ± 0.0073 | trestbps |
| 0.0166 ± 0.0146 | slope |
| 0.0156 ± 0.0039 | chol |
| 0.0107 ± 0.0073 | restecg |
| 0.0020 ± 0.0048 | fbs |
| 0.0010 ± 0.0189 | age |

Fig 5.1.2.4 AdaBoost using eli5

The above Fig shows that ca, thal, cp, oldpeak and thalach are top 5 important features. The weight of "ca" is 0.0966 ± 0.0357. The model performance is decreased with a random shuffling is 0.0966 and the performance varied from one-reshuffling to the next is 0.0357.

**5.1.2.5 Multilayer Perceptron using eli5:**

| Weight | Feature |
|---|---|
| 0.0683 ± 0.0269 | sex |
| 0.0576 ± 0.0455 | ca |
| 0.0576 ± 0.0434 | cp |
| 0.0312 ± 0.0117 | restecg |
| 0.0273 ± 0.0358 | slope |
| 0.0263 ± 0.0287 | oldpeak |
| 0.0254 ± 0.0292 | thal |
| 0.0166 ± 0.0078 | fbs |
| 0.0166 ± 0.0159 | age |
| 0.0137 ± 0.0305 | thalach |
| 0.0137 ± 0.0168 | trestbps |
| 0.0107 ± 0.0129 | exang |
| 0.0029 ± 0.0132 | chol |

Fig 5.1.2.5 : Multilayer Perceptron using eli5

The above Fig shows that sex, ca, cp, restecg and slope are top 5 important features. The weight of "sex" is 0.0683 ± 0.0269. The model performance is decreased with a random shuffling is 0.0683 and the performance varied from one-reshuffling to the next is 0.0269.
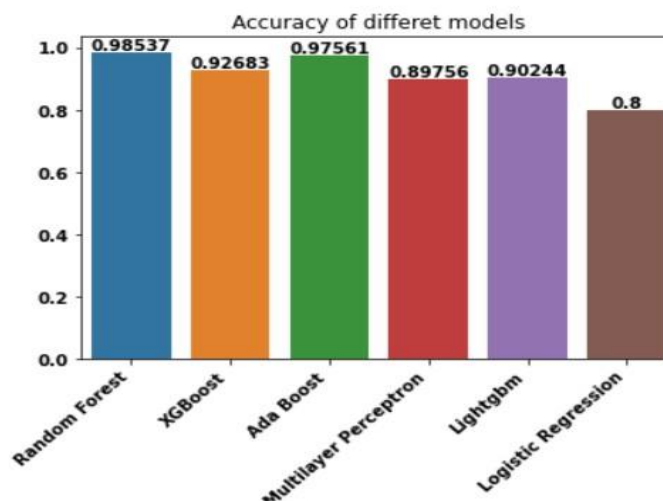


Chart-6.1: Accuracy of different models

The above graph shows the Comparison of Accuracy between Random Forest[0.98537], XGBoost [0.92683], AdaBoost [0.97561], Multilayer Perceptron [0.89756], LightGBM [0.90244] and Logistic Regression [0.8].
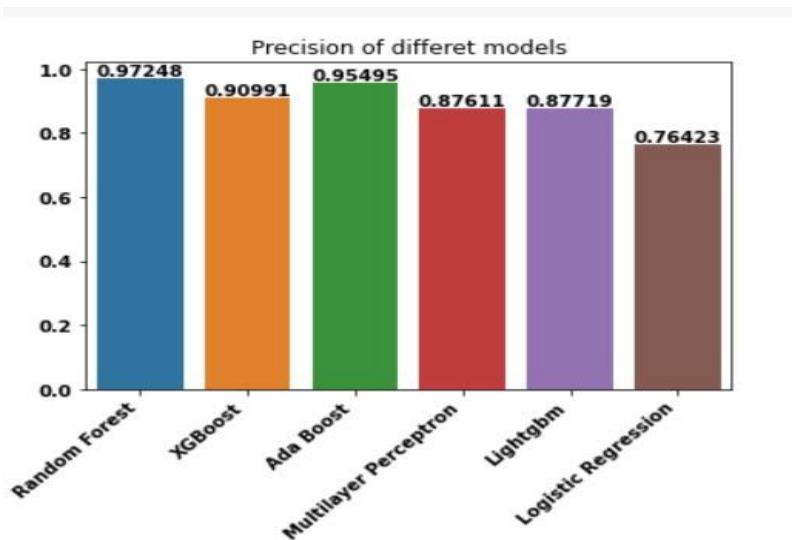
Chart-6.2: Precision of different models

The above graph shows the Comparison of precision between Random Forest[0.97248], XGBoost [0.90991], AdaBoost [0.95495], Multilayer Perceptron [0.87611], LightGBM [0.87719] and Logistic Regression [0.76423].
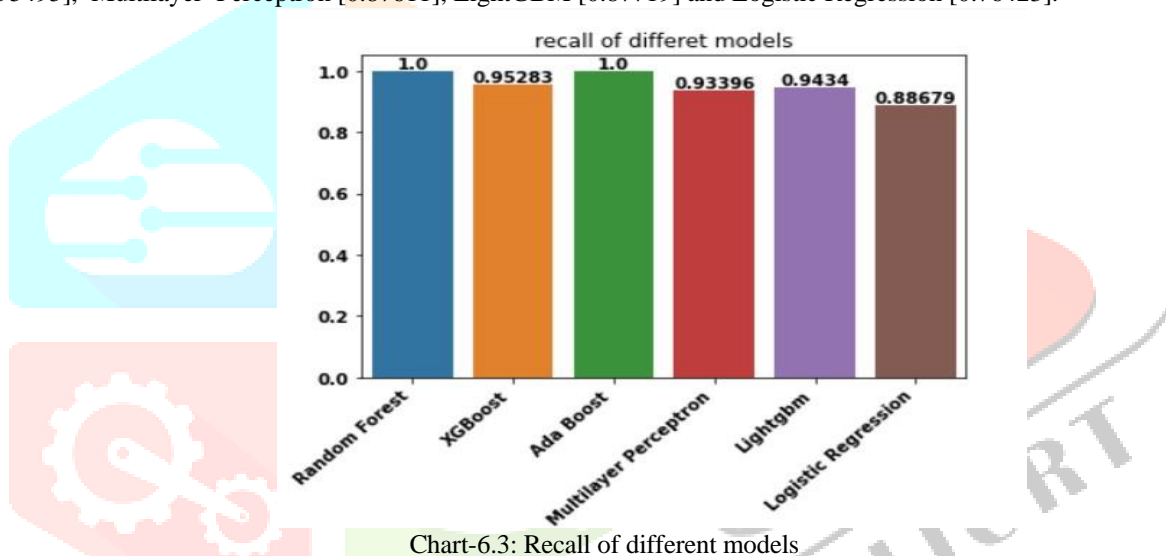


Chart-6.3: Recall of different models

The above graph shows the Comparison of recall between Random Forest [1.0],XGBoost [0.95283], AdaBoost [1.0], Multilayer Perceptron [0.93396], LightGBM [0.9434] and Logistic Regression [0.88679].
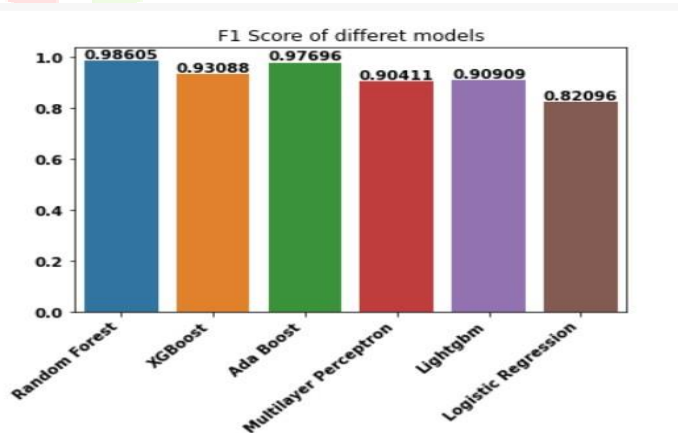


Chart-6.4: F1Score of different models

The above graph shows the Comparison of f1 score between Random Forest[0.98605], XGBoost [0.93088], AdaBoost [0.97696], Multilayer Perceptron[0.90411], LightGBM [0.90909] and Logistic Regression [0.82096].

# 7. CONCLUSION

In this paper, we have proposed the interpretation of machine learning algorithms using translation models. Model Interpretation not only helps to correct mistakes in your model and make your life as a machine learning engineer easier but also helps to build trust between people and the most important model as machine learning is used in a growing number of industries. This model relies on a heart rate database that includes patient data, age, gender, chol, treetops, and more.

We have used five machine learning algorithms namely Random Forest Classifier, AdaBoost, XGBoost, Light GBM and Multilayer Perceptron. For each algorithm we used two Interactive models namely Local Interpretable Model-agnostic Explanations (LIME) and Permission Value (Eli5) to make it descriptive.

In this Paper, we have used four key Performance Metrics namely Accuracy, Accuracy, Recall and F1 Score. Strategies for Random Forest Planning Operations are better available than other classification strategies. Random Jungle accuracy is 0.98537 higher than all other Algorithms.

We used Lime and eli5 in a random forest, after using LIME we got Prediction healthy chances are 0.05 and disease is 0.95. After inserting eli5 we obtained $0.0985 \pm 0.0285$. The performance of the model is reduced by a random shift of 0.0985 and the performance is divided from one reversal to the next by 0.0285. In, Our system Shows The Random Forest Algorithm Is Very Ideal Even In Difficult Times To Predict Heart Disease.

# 7. REFERENCES

[1].Aamodt, Agnar, and Enric Plaza. "Case-based reasoning: Foundationalissues, methodologicalvariations, and system approaches." AI communications7.1 (1994): 39-59.

[2].Adebayo, Julius, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. "Sanity checks for saliency maps." arXiv preprint arXiv:1810.03292 (2018).

[3].Alain, Guillaume, and Yoshua Bengio. "Understanding intermediate layers using linear classifier probes." arXiv preprint arXiv:1610.01644 (2016).

[4].Alber, Maximilian, Sebastian Lapuschkin, Philipp Seegerer, Miriam Hägele,Kristof T. Schütt,Grégoire Montavon, Wojciech Samek, Klaus-Robert Müller, Sven Dähne, and Pieter- Jan Kindermans. "iNNvestigate neural networks!." J. Mach. Learn. Res. 20, no. 93 (2019):1-8.

[5].Alberto, Túlio C, Johannes V Lochter, and Tiago A Almeida. "Tubespam: comment spam filtering on YouTube." In Machine Learning and Applications (Icmla), Ieee 14th International Conference on, 138–43. IEEE. (2015).

[6].Alvarez-Melis, David, and Tommi S. Jaakkola. "On the robustness of interpretability methods."arXiv preprint arXiv:1806.08049 (2018).

[7].Ancona, Marco, et al. "Towards better understanding of gradient-based attribution methods fordeep neural networks." arXiv preprint arXiv:1711.06104(2017).

[8].Apley, Daniel W., and Jingyu Zhu. "Visualizing the effects of predictor variables in black box supervised learning models." Journal of the Royal Statistical Society: Series B (Statistical Methodology) 82.4 (2020): 1059-1086.

[9].Athalye, Anish, and Ilya Sutskever. "Synthesizing robust adversarial examples." arXiv preprintarXiv:1707.07397 (2017).

[10].Breiman, Leo."Random Forests." Machine Learning 45 (1). Springer: 5 - 32 (2001).

[11]. Dandl, Susanne, Christoph Molnar, Martin Binder, Bernd Bischl. "Multi-objective counterfactual explanations". In: Bäck T. et al. (eds) Parallel Problem Solving from Nature – PPSN XVI. PPSN 2020. Lecture Notes in Computer Science, vol 12269. Springer, Cham (2020).

[12]. Deb, Kalyanmoy, Amrit Pratap, Sameer Agarwal and T. Meyarivan, "A fast and elitist multiobjective genetic algorithm: NSGA-II," in IEEE Transactions on Evolutionary Computation, vol. 6, no. 2, pp. 182-197, (2002).

[13]. Fernandes, Kelwin, Jaime S Cardoso, and Jessica Fernandes. "Transfer learning with partial observability applied to cervical cancer screening." In Iberian Conference on Pattern Recognition and Image Analysis, 243–50. Springer. (2017).

[14]. Fisher, Aaron, Cynthia Rudin, and Francesca Dominici. "All models are wrong, but many are useful: Learning a variable's importance by studying an entire class of prediction models simultaneously." (2018).

[15]. Goldstein, Alex, Adam Kapelner, Justin Bleich, and Emil Pitkin. "Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation." journal of Computational and Graphical Statistics 24, no. 1 (2015): 44-65.