



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 14 **Issue:** VI **Month of publication:** June 2026

DOI: <https://doi.org/10.22214/ijraset.2026.83585>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Consistency-Aware Evaluation of LIME and SHAP Explanations Using a Hybrid Similarity Metric

Anushka Arora, Hirday Pratap Singh, Mr. Avaneesh Kumar

Department of Computer Science and Engineering Galgotias University, Greater Noida, India

Abstract—*Explainable Artificial Intelligence methods like LIME and SHAP are used a lot to understand what machine learning models predict.. Most research focuses on generating explanations not checking if they are consistent across different methods. This paper presents a framework called Consistency-Aware Explainable AI that checks how well different explanation methods agree. It uses a score called Explanation Consistency Score, which combines two metrics.*

We tested this framework on the UCI Heart Disease dataset using Logistic Regression, Random Forest and XGBoost classifiers.

We evaluated explanation consistency across 50 test instances for each model.

The results show that being good at predicting does not mean an explanation is consistent. Logistic Regression got the Explanation Consistency Score while Random Forest got the lowest even though it was best at classifying.

The Explanation Consistency Score provides a reliable way to check explanation consistency across machine learning models.

The proposed ECS framework provides a computationally lightweight approach, for evaluating explanation consistency across machine learning models.

Index Terms—*Explainable AI, LIME, SHAP, Explanation Consistency, Feature Importance, Jaccard Similarity, Spearman Rank Correlation, Machine Learning*

I. INTRODUCTION

Machine learning models are being used more and more in areas like healthcare and finance. Some of these models like the ones that use different methods together or the ones that are really complex are hard to understand. This makes it tough for the people using these models to figure out why they are making predictions with the machine learning models. The machine learning models are like a box that you cannot see inside so it is hard to trust the predictions, from the machine learning models. [3].

Explainable Artificial Intelligence is a way to make Artificial Intelligence understandable. It does this by showing how the Artificial Intelligence makes its predictions based on the information it gets. There are a couple of ways that people usually do this. One way is called LIME (Local Interpretable Model-Agnostic Explanations) [1]. Another way is called SHAP (SHapley Additive exPlanations) [2]. Even though both LIME and SHAP try to explain how the Artificial Intelligence works they are very different. LIME and SHAP are different because they think about the problem in ways and use different methods to solve it. LIME makes a model just for one prediction at a time. On the hand SHAP uses a concept, from game theory to figure out how much each piece of information contributes to the prediction.

A big question is how often these two methods give the results. We need to know if LIME and SHAP agree with each other.. We have to see if LIME and SHAP agree in different situations like with different model architectures. If LIME and SHAP give different explanations, for the same prediction then it is hard to know which method to use when making decisions.

Contributions- This paper makes the following contributions:

- We introduce the Explanation Consistency Score (ECS) which is a hybrid metric combining Jaccard Similarity and Spearman Rank Correlation to measure the overlap between the feature sets and ranking agreement between LIME and SHAP explanations.
- We design a CA-XAI evaluation framework that works with any model and can be used to compare any two methods that explain which features are important. The CA-XAI evaluation framework is really flexible. Can be applied to any pair of feature-importance-based explanation methods.
- We are looking at how LIME and SHAP work with Logistic Regression, Random Forest and XGBoost on the UCI Heart Disease dataset. We want to see how consistent they are when we look at each instance and the whole model.
- We found out that when we use complex models to make predictions it can be hard to get consistent explanations. This means that there is a balance, between making predictions and getting explanations that make sense for LIME and SHAP.

The rest of the paper is organized as follows. Section II discusses the research gap. Section III reviews related work. Section IV describes the proposed methodology. Section V details the experimental setup. Section VI presents results and discussion. Section VII concludes the paper.

II. RESEARCH GAP AND MOTIVATION

Despite the growing use of techniques like LIME and SHAP, most research focuses on creating explanations without checking if they're reliable or consistent across different methods. Current methods check how good explanations are, how strong they are, or how easy to understand they are, but rarely see if different explanation methods give the insights for the same prediction. This is a problem in areas like healthcare and finance where wrong explanations can lead to bad decisions. Most studies do not have a way to measure if different explanation methods give similar results.

To fix this, we propose a metric called the Explanation Consistency Score (ECS). ECS uses Jaccard Similarity and Spearman Rank Correlation to check if LIME and SHAP explanations are consistent across different machine learning models.

III. RELATED WORK

A. Explanation Methods

Ribeiro and other people introduced LIME [1]. LIME is a way to understand how a complex model makes decisions. It does this by using a model that is easier to understand. This simpler model is usually a model with only a few important features. People like LIME because it is easy to use and it works with different types of models.

Lundberg and Lee [2] introduced SHAP. SHAP is a way to figure out which features are important for a model's decisions. It uses ideas from game theory to do this. SHAP is fair and consistent because it follows some rules. These rules are called efficiency, symmetry, and additivity. For models that use trees, there is an algorithm called TreeSHAP. TreeSHAP makes it possible to calculate these values exactly and quickly [11].

Some people made a list of all the ways to explain how models work. Adadi and Berrada [3] did this. They grouped these methods by what they do, how they work with models, and what they produce. Guidotti et al. [4] and other people also made a list of ways to explain models. They talked about the challenges of figuring out if these methods are working correctly. They discussed ways to explain models and the problems that come with it.

B. Evaluation of Explanation Quality

Evaluating the quality of explanations is still something we are trying to figure out. We have some ways to measure explanations, like how they match the truth, if they are stable, and if they are easy to understand. We do not have one standard way to judge all explanations. The problem is that different methods to explain things can show parts of how a model works, so it is hard to compare them directly without a clear way to measure them. Evaluating the quality of explanations is really important. We need a better way to do it.

C. Consistency Between Explanation Methods

The problem of getting explanation methods to agree with each other is something that people haven't looked into very much. When we use methods like LIME and SHAP on the same model, we find that they often pick similar but not exactly the same important features. However, nobody has come up with a way to measure how well these methods agree with each other. This work is trying to fix this issue by introducing something called ECS, which is a simple and easy-to-understand way to measure how consistent these methods are, and this thing is called a consistency metric for explanation methods, like LIME and SHAP.

IV. PROPOSED METHODOLOGY

Fig. 1 shows the process of the proposed CA-XAI framework. The CA-XAI framework has four parts. First, there is the data preprocessing stage. Then, the CA-XAI framework goes through the model training stage, and after that, it does the explanation generation stage. Finally, the CA-XAI framework does the ECS computation stage.

A. Dataset

The UCI Heart Disease Dataset [12] is used for experiments. This dataset has 303 records. It includes 13 features, for each patient.

These features are the patients age the patients sex, the type of chest pain the patient has the patients resting blood pressure the patients serum cholesterol level the patients fasting blood sugar level the results of the patients resting ECG the patients heart rate whether the patient gets exercise-induced angina the patients ST depression the patients ST slope the number of major vessels the patient has and the patients thalassemia status. The UCI Heart Disease Dataset has a target variable. This binary target variable shows whether the patient has heart disease or not.

B. Classification Models

Three classifiers with varying complexity are evaluated:

- Logistic Regression (LR) [7]: A linear classifier used as the baseline. Its linear structure makes LIME's local surrogate and SHAP's additive attributions theoretically compatible.
- Random Forest (RF) [5]: A non-linear ensemble of decision trees using bagging. Diversity among trees introduces complexity in local explanation approximation.
- XGBoost (XGB) [6]: A gradient-boosted tree ensemble known for high predictive performance. TreeSHAP [11] is used for efficient Shapley value computation.

C. Explainability Methods

1) **LIME**: LIME generates instance-level explanations by perturbing an input sample, querying the black-box model on the perturbed samples, and fitting a locally weighted interpretable surrogate.

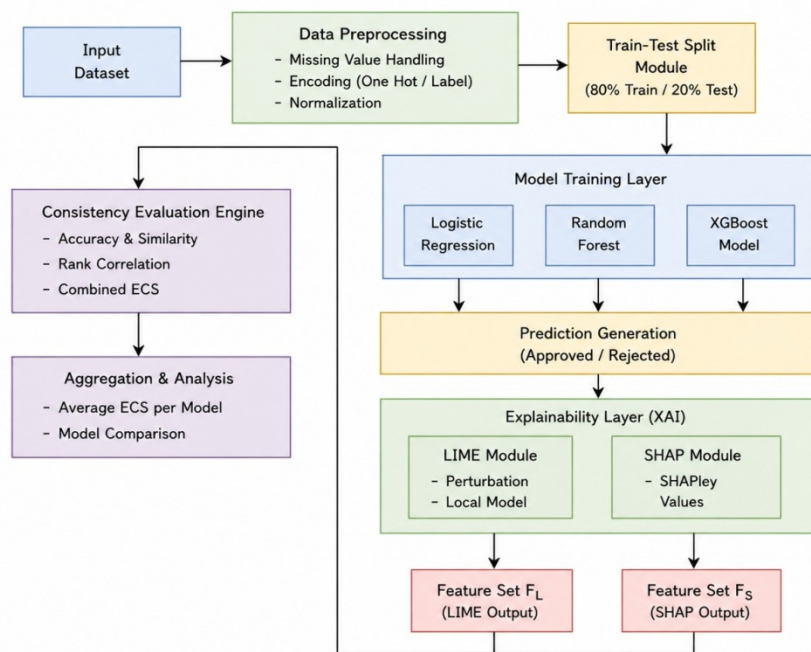


Fig. 1. Overall pipeline of the proposed CA

XAI framework including preprocessing, model training, LIME/SHAP explanation generation, and ECS computation.

LIME internally generates perturbed samples around a target input instance and observes the prediction behavior of the black-box model on these samples. Similarity weights are assigned based on the distance between perturbed samples and the original instance. A locally interpretable surrogate model is then fitted to approximate model behavior in the local neighborhood.

The optimization objective of LIME is:

$$\arg \min L(f, g, \pi_x) + \Omega(g) \quad (1)$$

$$g \in G$$

where f is the black-box model, g is the surrogate model, π_x is the locality-aware weighting function, and $\Omega(g)$ penalizes surrogate complexity. The top- k features by coefficient magnitude are extracted from the fitted surrogate.

2) *SHAP*: *SHAP* computes feature importance using Shap-ley values derived from cooperative game theory. Each feature contribution is calculated as the average marginal contribution of that feature across all possible subsets of input features. The Shapley value for feature i is:

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} [f(S \cup \{i\}) - f(S)] \quad (2)$$

where F is the complete feature set, S is a subset of features, and $f(S)$ is the model output using only features in S . For tree-based models, the TreeSHAP algorithm [11] used for efficient computation.

D. Consistency-Aware XAI Algorithm

The workflow of the proposed CA-XAI framework is as follows.

Input: Dataset D , machine learning models M , top- k feature count k

Output: ECS for each model

- 1) Load and preprocess dataset.
- 2) Encode categorical features; normalize numerical features.
- 3) Split dataset into training and testing sets.
- 4) Train models: Logistic Regression, Random Forest, XG-Boost.
- 5) For each test instance x_i :
 - Generate LIME explanation; extract top- k features F_L .
 - Generate SHAP explanation; extract top- k features F_S .
 - Compute Jaccard Similarity J .
 - Compute Spearman Rank Correlation ρ .
 - Normalize ρ to $[0, 1]$.
 - Compute instance-level ECS.
- 6) Aggregate ECS scores across all instances.
- 7) Report mean ECS and standard deviation per model.

E. Explanation Consistency Score (ECS)

The ECS captures two parts of explanation agreement.

a) *Jaccard Similarity*: Let F_L be the top- k features from LIME and F_S from SHAP. The Jaccard Similarity measures the feature-set overlap [8]:

$$J = \frac{|F_L \cap F_S|}{|F_L \cup F_S|} \quad (3)$$

$J \in [0, 1]$, where 1 indicates identical feature sets.

b) *Spearman Rank Correlation*: For features shared by both F_L and F_S , the Spearman Rank Correlation measures ranking consistency [9]:

$$\rho = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)} \quad (4)$$

where d_i is the rank difference for shared feature i and n

is the number of shared features. $\rho \in [-1, 1]$.

c) *ECS Formulation*: The Spearman Rank Correlation value is linearly normalized from the range $[-1, 1]$ to $[0, 1]$ before combining it with Jaccard Similarity:

$$\rho_{\text{norm}} = \frac{\rho + 1}{2} \tag{5}$$

The final ECS is computed as:

$$\text{ECS} = 0.5 \cdot J + 0.5 \cdot \rho_{\text{norm}} \tag{6}$$

$\text{ECS} \in [0, 1]$, in which the ECS is higher it means that the two explanation methods are more consistent with each other. The ECS gives weight to things, which means it thinks that the identity of the features and the order they are, in are equally important.

F. Computational Complexity

For top- k features, Jaccard Similarity needs $O(k)$ operations. Spearman Rank Correlation requires $O(k \log k)$ operations because it has to rank features. The overall ECS computation complexity, per explanation pair is therefore $O(k \log k)$ which makes it lightweight relative to the cost of generating LIME and SHAP explanations.

V. EXPERIMENTAL SETUP

A. Preprocessing

We used one-encoding for the categorical variables. This is a way to make them work with our model. The continuous features were also changed. We made sure they had a mean of zero and a variance of one. This is called z-score normalization. We did this before we trained our model. We did it to make sure the continuous features were on the scale, as the categorical variables.

B. Train-Test Split

We took the dataset. Split it into two parts. We used 80 percent for training and 20 percent for testing. This way the class distribution stays the same. When we were picking the settings for the model we used a special method, on the training part. This method is called 5-fold cross-validation.

C. Explanation Generation

LIME and SHAP explanations were made for fifty test instances that were chosen randomly for each model. For each test instance the top ten features that had the biggest impact were picked from each method. This was done by looking at how each feature contributed to the result. The ECS was calculated for each test instance. Then combined together the average and standard deviation for each model. The performance of the classification was checked using measures, like how accurate it was how precise it was how well it recalled things and the F1-score of LIME and SHAP explanations.

D. Implementation Details

The proposed framework was built using Python. I used Google Colab with Python version 3.11. Here are the libraries I used:

- scikit-learn[10] is used for preprocessing and model training
- shap is used for feature attribution generation
- lime is used for local explanation generation
- xgboost is used for gradient-boosted classification
- scipy is used for Spearman rank correlation computation

TA used a fixed seed of 42 everywhere to make sure the results can be repeated. The following hyperparameters were used:

- LogisticRegression: max_iter=1000
- RandomForest: default scikit-learn parameters
- XGBoost: eval_metric=logloss

VI. RESULTS AND DISCUSSION

A. Classification Performance and ECS

Table I shows how well the models did with classification. It also gives us the ECS values for each of the three models.

TABLE I CLASSIFICATION PERFORMANCE AND EXPLANATION CONSISTENCY ACROSS MODELS

Model	Acc.	Prec.	Rec.	F1	MeanECS	StdECS
LogisticReg	0.853	0.900	0.841	0.869	0.5576	0.0930
RandomFor	0.891	0.922	0.887	0.904	0.3879	0.1744
XGBoost	0.875	0.903	0.878	0.891	0.5125	0.1289

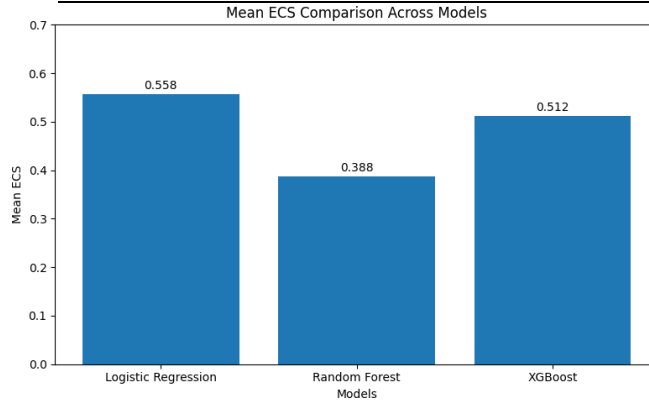


Fig. 2. Mean ECS comparison across Logistic Regression, Random Forest, and XGBoost models. Logistic Regression achieved the highest explanation consistency, while Random Forest exhibited the lowest ECS despite higher predictive accuracy.

B. Analysis

The results of the experiment show that the explanations given by machine learning models are not consistent. Logistic Regression gave the consistent explanations with a score of 0.5576. This means that the explanations from LIME and SHAP mostly agreed with each other for Logistic Regression. XGBoost had a score of 0.5125 which's pretty good.. Random Forest had the lowest score of 0.3879 which means its explanations were not very consistent.

What is interesting is that Random Forest was really good at making predictions with an accuracy of 89.13 percent.. It had the lowest explanation consistency score. This shows that just because a model is good at making predictions it does not mean that its explanations will always make sense. Models like Random Forest that combine different models can make it hard for LIME and SHAP to agree on which features are important. This is because these models have rules, for making decisions.

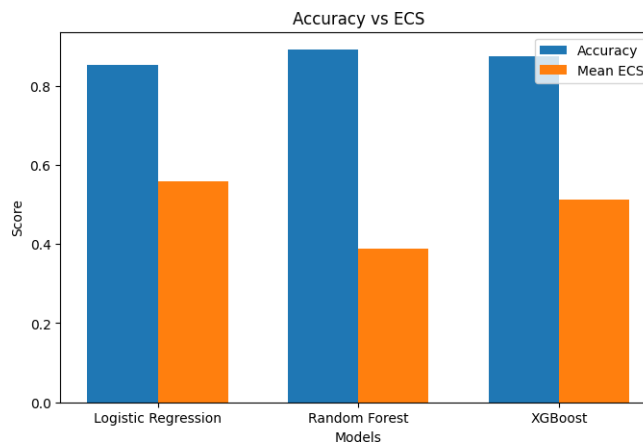


Fig. 3. Comparison between classification accuracy and mean ECS across Logistic Regression, Random Forest, and XGBoost models

In contrast Logistic Regression gave predictions but had consistent explanations. This is because it uses an clear method.

Both LIME and SHAP should give results for simple modelslikeLogisticRegression.LIMEmakesasimplemodel to mimic the one and SHAP adds up the effects of each feature.Theseresultsshowthatyoucan'thavebothpredictions and clear explanations at the same time. This is important when choosing a model, for applications where understanding the model matters.

C. ExplanationStability

The scores from ECS give us an idea of how stable the explanationsre when we test them. Logistic Regression hadthe amount of variation in ECS scores, which is 0.0930. This meansthattheexplanationsfromLogisticRegressionaremore stable when we look at samples. On the hand Random Forest hadthebiggestamountofvariationinECSscores,which is0.1744.

This suggests that the explanations from Random Forest are not as consistent and can change a lot when we use models that combine many things. XGBoost had an amount of variation, in ECS scores, which is 0.1289.

The heart disease dataset is something that the classifiersre really good at predicting.If you look at the confusion matrices in Fig. 4you can see that all three classifiers do a job. The RandomForestheartdiseaseclassifieristhebest,atgettingthe positive instances right.The Logistic Regression heart disease classifier is also good because it makes predictions and does not make a lot of unstable classifications. The heart disease dataset is predicted well by all three classifiers.

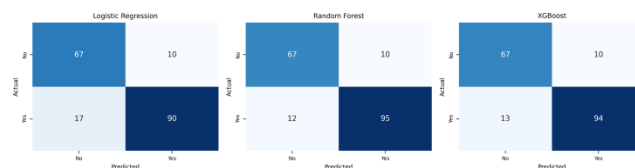


Fig.4.ConfusionmatricesforLogisticRegression,RandomForest,andXGBoost classifiers evaluated on the UCI Heart Disease dataset.

Fig. 5shows how ECS values are spread out across the test instances we sampled. Logistic Regression has a distribution that is more grouped with less variation which means that LIMEandSHAPexplanationsagreewitheachotherinaway. Random Forest shows wider dispersion which suggests less stable explanation consistency across all samples.

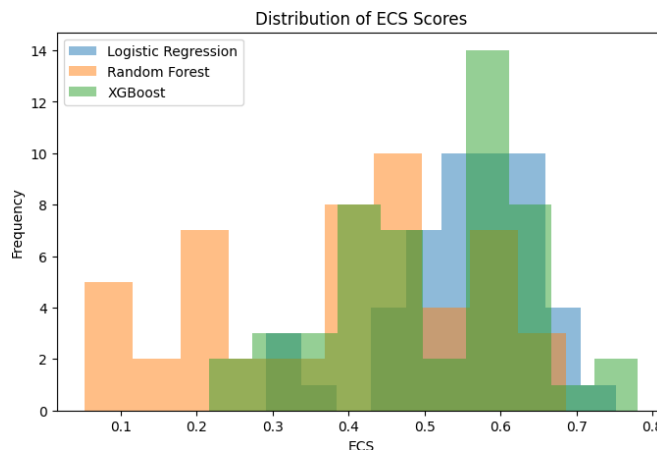


Fig.5.Distributionofper-instanceECSscoresacross50sampledtestinstances for Logistic Regression, Random Forest, and XGBoost models.

Fig. 6shows how much the explanations from LIME and SHAP agree on each feature. If we look at things like the type of chest pain the ST slope and if someone gets angina when they exercise we can see that LIME and SHAP explanations usually agree on these things. This means that these features areimportantforthemodelsandtheyareimportanteverytime. The features like chest pain type and ST slope are consistent. This is true, for LIME explanations and SHAP explanations.

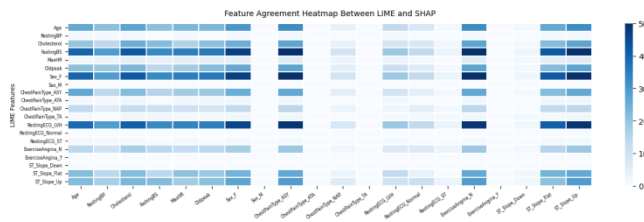


Fig.6.

VII. CONCLUSION AND FUTURE WORK

This paper is about the Explanation Consistency Score, which's a way to measure how well two explanation methods, LIME and SHAP agree with each other. The Explanation Consistency Score is a mix of two things: Jaccard Similarity and Spearman Rank Correlation.

The researcher tried out the Explanation Consistency Score with a few models, like Logistic Regression, Random Forest and XGBoost on a dataset about heart disease. They found that Logistic Regression had the Explanation Consistency Score, which was 0.5576 and it was also very stable. On the hand Random Forest was really good at predicting things but its Explanation Consistency Score was the lowest, at 0.3879.

This makes us think that maybe models that are too complex do not do a job of explaining things in a consistent way. So when we are choosing a model we should not just think about how it predicts things but also about how well it explains things.

The good thing about the Explanation Consistency Score is that it is not hard to compute and it works with any model. It is also easy to add to the way we already evaluate models.

There are a few things to keep in mind though. The researchers only tried this out on one dataset so we do not know if it will work the way on other datasets. They also only used LIME and SHAP so we do not know what would happen with other explanation methods. They gave equal weight to Jaccard Similarity and Spearman Rank Correlation which might not always be the best thing to do.

The researchers did everything they could to make sure their results are reliable, by using the random seeds every time and making their code public.

In the future the researchers want to try the Explanation Consistency Score with complex models like deep neural networks and see if they can make it work better by changing the way they weigh Jaccard Similarity and Spearman Rank Correlation.

REFERENCES

- [1] M. T. Ribeiro, S. Singh, and C. Guestrin, "Why Should I Trust You?: Explaining the Predictions of Any Classifier," in Proc. ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining, pp. 1135–1144, 2016.
- [2] S. M. Lundberg and S.-I. Lee, "A Unified Approach to Interpreting Model Predictions," in Advances in Neural Information Processing Systems (NIPS), vol. 30, pp. 4765–4774, 2017.
- [3] A. Adadi and M. Berrada, "Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI)," IEEE Access, vol. 6, pp. 52138–52160, 2018.
- [4] R. Guidotti, A. Monreale, S. Ruggieri, F. Turini, F. Giannotti, and Pedreschi, "A Survey of Methods for Explaining Black Box Models," ACM Computing Surveys, vol. 51, no. 5, pp. 1–42, 2019.
- [5] L. Breiman, "Random Forests," Machine Learning, vol. 45, no. 1, pp. 5–32, 2001.
- [6] T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," in Proc. ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining, pp. 785–794, 2016.
- [7] D. W. Hosmer and S. Lemeshow, Applied Logistic Regression, 2nd ed. New York, NY: Wiley, 2000.
- [8] P. Jaccard, "The Distribution of Flora in the Alpine Zone," New Phytologist, vol. 11, no. 2, pp. 37–50, 1912.
- [9] C. Spearman, "The Proof and Measurement of Association Between Two Things," American Journal of Psychology, vol. 15, no. 1, pp. 72–101, 1904.
- [10] F. Pedregosa et al., "Scikit-learn: Machine Learning in Python," Journal of Machine Learning Research, vol. 12, pp. 2825–2830, 2011.
- [11] S. M. Lundberg et al., "From Local Explanations to Global Understanding with Explainable AI for Trees," Nature Machine Intelligence, vol. 2, no. 1, pp. 56–67, 2020.
- [12] A. Janosi, W. Steinbrunn, M. Pfisterer, and R. Detrano, "Heart Disease Dataset," UC Machine Learning Repository, 1988. [Online]. Available: <https://archive.ics.uci.edu/ml/datasets/heart+disease>



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)