

Article

Not peer-reviewed version

Enhancing SQL Injection Detection with Trustworthy Ensemble Learning and Boosting Models Using Local Explanation Techniques

[Thi-Thu-Huong Le](#)*, [Yeonjeong Hwang](#), Changwoo Choi, Rini Wisnu Wardhani, [Dedy Septono Catur Putranto](#), [Howon Kim](#)*

Posted Date: 24 October 2024

doi: 10.20944/preprints202410.1878.v1

Keywords: Explained AI; SQL Injection Detection; Decision Tree; Random Forest; XGBoost; AdaBoost; Gradient Boosting Decision Tree; Histogram Gradient Boosting Decision Tree; Local Explanation; SHAP; LIME



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

Enhancing SQL Injection Detection with Trustworthy Ensemble Learning and Boosting Models Using Local Explanation Techniques [†]

Thi-Thu-Huong Le ^{1,*} , Yeonjeong Hwang ² , Changwoo Choi ^{2,3} , Rini Wisnu Wardhani ² , Dedy Septono Catur Putranto ⁴  and Howon Kim ^{2,*} 

¹ Blockchain Platform Research Center, Pusan National University, Busan 609735, South Korea; lehuong7885@gmail.com

² School of Computer Science and Engineering, Pusan National University, Busan 609735, South Korea; yeonjeong@islab.re.kr; changwoo7463@gmail.com; rini.wisnu@pusan.ac.kr

³ SmartM2M, Busan 48058, South Korea; changwoo7463@gmail.com;

⁴ IoT Research Center, Pusan National University, Busan 609735, South Korea; dedy.septono@pusan.ac.kr;

* Correspondence: huongle@pusan.ac.kr (Le, T.T.H.); howonkim@pusan.ac.kr (Kim, H.)

[†] This paper is an extended version of our paper published In Proceedings of the 2024 International Conference on Platform Technology and Service (PlatCon-24), Jeju, Korea, 26–28 August 2024.

Abstract: This paper presents a comparative analysis of several decision models for detecting Structured Query Language (SQL) injection attacks, which remain one of the most prevalent and serious security threats to web applications. SQL injection enables attackers to exploit databases, gain unauthorized access, and manipulate data. Traditional detection methods often struggle due to the constantly evolving nature of these attacks, the increasing complexity of modern web applications, and the lack of transparency in the decision-making processes of machine learning models. To address these challenges, we evaluated the performance of various models, including Decision Tree, Random Forest, XGBoost, AdaBoost, Gradient Boosting Decision Tree (GBDT), and Histogram Gradient Boosting Decision Tree (HGBDT), using a comprehensive SQL injection dataset. The primary motivation behind our approach is to leverage the strengths of ensemble learning and boosting techniques to enhance detection accuracy and robustness against SQL injection attacks. By systematically comparing these models, we aim to identify the most effective algorithms for SQL injection detection systems. Our experiments show that Decision Tree, Random Forest, and AdaBoost achieved the highest performance, with an accuracy of 99.50% and an F1 score of 99.33%. Additionally, we applied SHapley Additive exPlanations (SHAP) and Local Interpretable Model-agnostic Explanations (LIME) for local explainability, illustrating how each model classifies normal and attack cases. This transparency enhances the trustworthiness of our approach to detecting SQL injection attacks. These findings highlight the potential of ensemble methods to provide reliable and efficient solutions for detecting SQL injection attacks, thereby improving the security of web applications.

Keywords: explained AI; SQL injection detection; decision tree; random forest; XGBoost; AdaBoost; Gradient Boosting Decision Tree; Histogram Gradient Boosting Decision Tree; local explanation; SHAP; LIME

1. Introduction

Structured Query Language injection (SQLi) remains one of the most critical security vulnerabilities affecting web applications, consistently ranking among the top security threats in reports like the OWASP Top Ten [1]. SQL injection exploits weaknesses in input validation, enabling attackers to manipulate SQL queries to access, modify, or even delete sensitive data stored in databases [2,3]. The widespread use and serious consequences of SQL injection attacks make it crucial to develop robust detection and prevention mechanisms to safeguard web applications and their underlying databases.

Traditional approaches to mitigating SQL injection vulnerabilities, such as input validation, parameterized queries, and prepared statements, aim to sanitize user inputs and block malicious SQL commands from being executed [4,6,7]. However, these methods often fail in dynamically generated query environments and require rigorous implementation practices that are not always followed consistently [8]. This gap becomes more problematic as SQL injection techniques evolve and grow in

sophistication, underscoring the need for more advanced and adaptable detection mechanisms capable of handling emerging threats.

Machine learning (ML) has emerged as a promising solution to improve SQL injection detection by learning historical attack patterns and generalizing to new threats [9]. Among the various ML techniques, ensemble learning and boosting models have shown considerable potential in this domain, as they combine multiple weak learners into a stronger and more accurate classifier [12]. These models not only improve detection accuracy but also increase robustness against sophisticated SQL injection attacks, making them suitable for dynamic and complex web environments.

Ensemble learning methods, such as Decision Trees and Random Forests, take advantage of the collective wisdom of multiple models to achieve better performance than individual classifiers [10,13]. Boosting algorithms, including AdaBoost, Gradient Boosting Decision Trees (GBDT), and Histogram Gradient Boosting Decision Trees (HGBDT), further enhance this capability by sequentially focusing on misclassified instances, thereby improving the model's precision and recall [11,14,15]. These approaches offer a nuanced understanding of attack patterns, enabling them to adapt to new and sophisticated SQL injection techniques.

However, there remains a significant gap in the application of explainable machine learning models for SQL injection detection. Although numerous studies have focused on explaining the results of the model in the context of SQL injection detection [16,17], much of the research on explainable AI (XAI) has focused on other security domains, such as intrusion detection systems (IDS) [10,15,18]. Techniques such as SHapley Additive exPlanations (SHAP) [19] and Local Interpretable Model-agnostic Explanations (LIME) [20] have been widely applied in those fields to clarify model decisions and enhance transparency.

Despite this progress in other domains, explainable AI models have not yet been sufficiently explored or adapted for SQL injection detection. Given the critical importance of SQL injection attacks and the complexity of machine learning models to detect them, extending these XAI techniques to this domain is essential. By incorporating SHAP and LIME into SQL injection detection, we can improve the trustworthiness and interpretability of the model's decision-making process, thereby increasing confidence in its real-world deployment.

In this study, we focus on evaluating the performance of several ensemble learning and boosting models for SQL injection detection, specifically Decision Tree, Random Forest, XGBoost, AdaBoost, GBDT, and HGBDT. These models were selected due to their proven success in other security domains and ability to efficiently handle large, complex datasets. Our objective is to enhance SQL injection detection systems by improving precision, recall, and overall robustness in identifying malicious SQL queries. To further increase trust and explain the proposed model decision, we incorporate well-known explainability techniques, such as SHAP and LIME, to provide transparent insights into the decision-making process of the models.

The primary contributions of this paper are as follows:

- We conduct a comparative analysis of six decision models, including Decision Tree, Random Forest, XGBoost, AdaBoost, GBDT, and HGBDT, using a comprehensive SQL injection dataset.
- We evaluated the performance of these models based on key metrics such as precision, precision, recall, and F1 score, identifying the most effective machine learning techniques for SQL injection detection.
- We apply SHAP and LIME to explain the decision-making processes of each model, improving transparency and trustworthiness in SQL injection detection.
- We provide insights into the strengths and limitations of ensemble learning and boosting models for practical deployment in real-world SQL detection systems.

The remainder of this paper is structured as follows. Section 2 reviews related work on SQL injection detection using machine learning techniques. Section 3 outlines the methodology and experimental setup used for the study. Section 4 presents the experimental results, followed by a

discussion in Section 5. Finally, Section 6 concludes the paper and suggests directions for future research.

2. Related Work

Over the past two decades, the detection and prevention of SQL injection (SQLi) attacks have been the focus of extensive research. This section reviews significant contributions to the field, particularly those employing machine learning (ML) techniques, to highlight the progress and challenges in SQLi detection.

2.1. Machine Learning Techniques for SQLi Detection

Initial efforts to mitigate SQL injection attacks primarily relied on static analysis and heuristic-based methods. These approaches involve code reviews and the implementation of best practices such as input validation, parameterized queries, and prepared statements [2,4]. While effective in reducing vulnerabilities, these methods are often labor-intensive and prone to human error, especially in complex or dynamically generated queries.

The advent of machine learning has provided a more dynamic and adaptive approach to SQLi detection. ML models can learn from historical attack data to identify and mitigate new and evolving threats [5]. Several studies have explored the use of various ML algorithms for SQLi detection with promising results. Support Vector Machines (SVM) and Neural Networks, for instance, have been employed to classify SQL queries as malicious or benign based on features extracted from query structures and patterns. Valeur et al. [21] proposed an anomaly-based intrusion detection system leveraging SVMs to detect SQL injection and other web attacks. Similarly, Gao et al. [22] used a neural network model trained on labeled datasets of SQL queries to achieve high detection accuracy.

In addition to these, Decision Trees and Random Forests have been widely used due to their interpretability and robustness. Xu et al. [23] demonstrated that decision tree classifiers could effectively distinguish between malicious and benign SQL queries by learning the characteristic patterns of SQL injection attacks. Random forests further enhance this capability by aggregating the predictions of multiple decision trees, improving accuracy and reducing overfitting [13].

Ensemble learning methods, which combine multiple weak learners to form a strong classifier, have shown significant potential in SQLi detection. AdaBoost and Gradient Boosting Decision Trees (GBDT) are notable examples of boosting algorithms successfully applied in this domain. These models focus on instances that are difficult to classify, iteratively improving the classifier's performance. Pan et al. [24] used AdaBoost to enhance SQLi attack detection accuracy by adapting to new attack patterns. Le et al. [25] demonstrated the effectiveness of GBDT in handling complex datasets with imbalanced class distributions, achieving high-precision detection of SQLi attacks.

XGBoost and Histogram-based Gradient Boosting Decision Trees (HGBDT) have further advanced the state-of-the-art in SQLi detection. These models optimize the boosting process by efficiently handling large datasets and improving generalization capabilities. Chen and Guestrin [26] introduced XGBoost, a popular algorithm for SQLi detection due to its scalability and performance. Ke et al. [27] extended this approach with HGBDT, offering improved computational efficiency and accuracy for high-dimensional data.

Several comparative studies have evaluated the performance of different ML models for SQLi detection. These studies typically compare accuracy, precision, recall, and F1 score metrics to identify the most effective models. Nguyen et al. [28] comprehensively evaluated various ML algorithms, including decision trees, SVMs, and neural networks, for detecting SQL injection attacks, highlighting the superior performance of ensemble learning and boosting models.

Despite significant advancements, several challenges remain in detecting SQL injection attacks using ML techniques. One major challenge is the imbalanced nature of datasets, where non-malicious queries vastly outnumber malicious ones. This imbalance can lead to biased models that favor the

majority class. Techniques such as oversampling, undersampling, and synthetic data generation have been proposed to address this issue [29].

Another challenge is the evolving nature of SQL injection attacks. Attackers continuously develop new techniques to bypass existing detection mechanisms, requiring models that can adapt to these changes [30]. Continuous learning and real-time monitoring systems are potential solutions to enhance the adaptability and effectiveness of ML-based detection systems [9].

2.2. Challenges with Explainable AI for SQL Injection Detection

While machine learning models have shown great promise in detecting SQLi attacks, the challenge of interpretability and trustworthiness remains a key barrier to their widespread adoption. SQL injection detection, especially in security-critical web applications, requires not only high-performance models but also clear explanations of how these models make decisions. Explainable AI (XAI) has emerged as a solution to this problem, offering tools to interpret and visualize machine learning models. However, applying XAI in the context of SQL injection detection brings specific challenges.

One significant challenge in applying explainable AI models like SHapley Additive exPlanations (SHAP) [19] and Local Interpretable Model-agnostic Explanations (LIME) [20] to SQLi detection lies in the inherent complexity and variability of SQLi attacks. SQLi attacks often involve obfuscated, multi-layered queries, making it difficult to isolate the specific features or patterns that contribute to a detection decision. This obfuscation can hinder the creation of meaningful explanations that security analysts can easily interpret.

Moreover, one of the main objectives of applying XAI in SQLi detection is to increase the trust in the predictions of the model. In real-world deployments, security teams need to understand not only which features contributed to a prediction but also why the model classified a specific SQL query as malicious. Current XAI methods provide local explanations, which provide insight for individual queries but may not explain broader patterns or attack vectors. For example, explaining why one SQL query was flagged as malicious doesn't always provide insights into other, potentially related attacks. This limitation underscores the challenge of developing XAI methods that not only provide accurate local explanations but also offer global insights into SQLi attack trends.

In summary, while traditional methods provide a foundational layer of security, integrating machine learning techniques, particularly ensemble learning and boosting models, offers a more robust and dynamic approach to SQLi detection. The current research landscape highlights the effectiveness of these models in identifying and mitigating SQL injection attacks, paving the way for further advancements in this critical area of cybersecurity. However, the application of XAI in this domain remains a challenge due to the complex nature of SQLi attacks, which requires future research to address interpretability issues and improve trust in AI-driven security systems.

3. Methodology

This section presents our methodology for detecting SQL injection using ensemble learning and boosting models, depicted in Figure 1. The process begins with data preprocessing, where we encode string data into a numerical format using a Label Encoder and split the dataset into training and testing sets. We train six models: Decision Tree, Random Forest, XGBoost, AdaBoost, Gradient Boosting Decision Tree (GBDT), and Histogram Gradient Boosting Decision Tree (HGBDT) on the training dataset. These models are then saved as pretrained for inference on the testing dataset. Finally, we evaluate the performance of our approach using key metrics such as confusion matrix, Receiver Operating Characteristic (ROC) curve, precision, recall, accuracy, and F1 score on both the training and testing datasets.

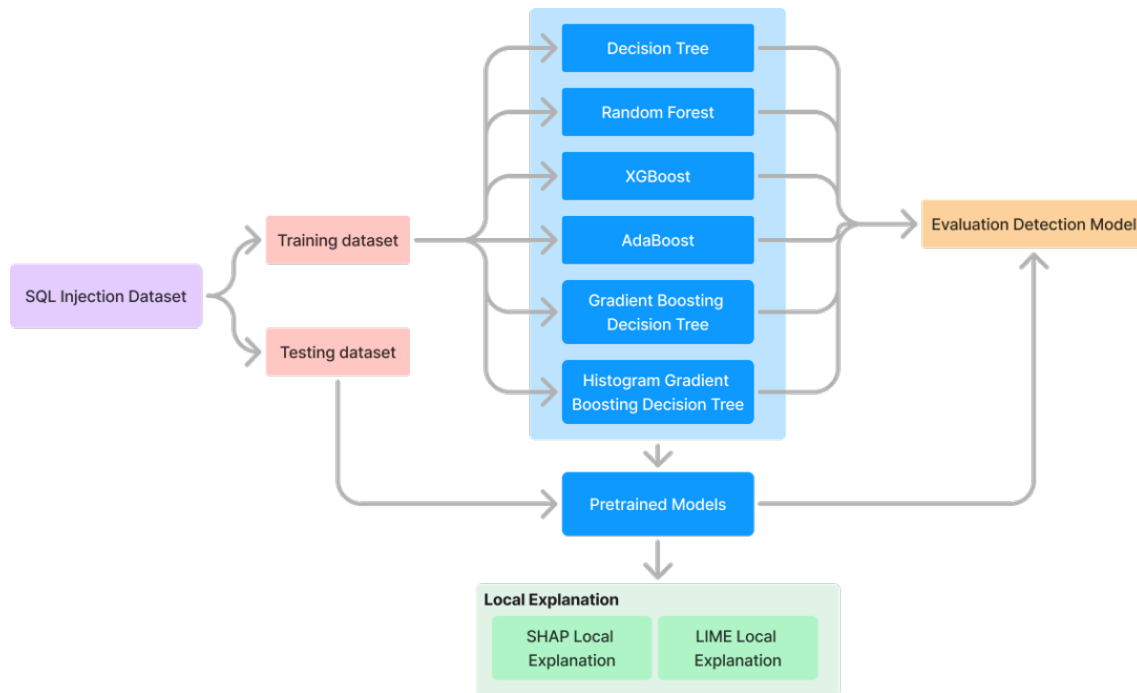


Figure 1. The proposed method for detecting SQL injection from ensemble learning and boosting models

3.1. Data Description

This dataset is sourced from the SQL Injection Attack Dataset [31]. It contains raw SQL query strings labeled with integers in the Label column, where a value of 0 indicates a non-malicious query, and a value of 1 denotes a malicious query. Figure 2 illustrates the distribution of sample data labeled as attack and non-attack in the SQL Injection dataset. The dataset comprises approximately 1,200 attack samples and around 20,000 non-attack samples. To evaluate the performance of our models, we partitioned the dataset into training and testing sets with a ratio of 70:30.

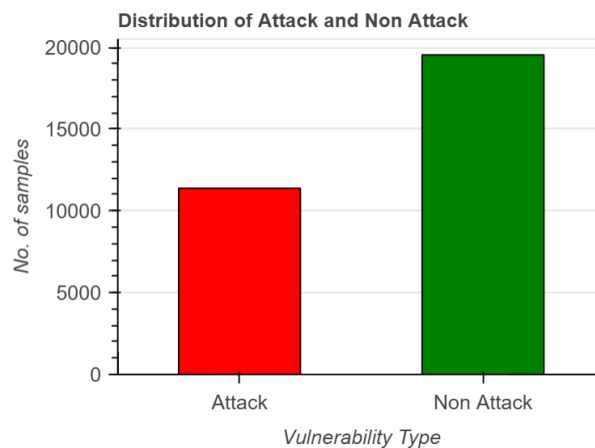


Figure 2. Distribution sample data.

3.2. SQL Vulnerabilities Detection Based Ensemble and Bagging Models

The rationale for selecting these models lies in their diverse strengths and complementary capabilities. Decision Trees provide straightforward interpretation and quick inference, making them useful for initial insights and rapid prototyping. Random Forests, an ensemble of Decision Trees, enhance performance by reducing overfitting and improving generalization. Boosting models like XGBoost and AdaBoost iteratively focus on misclassified instances, increasing the overall model accuracy and

robustness. GBDT and HGBDT extend boosting methods by efficiently handling large datasets and high-dimensional feature spaces, making them suitable for complex real-world applications.

- **Decision Tree.** This model has hierarchical structures where each node represents a feature, and each branch represents a decision rule based on that feature. The decision tree recursively partitions the feature space into regions that minimize a splitting criterion (e.g., Gini impurity or entropy).
- **Random Forest.** It is an ensemble learning method that constructs multiple decision trees during training and outputs the mode of the classes (classification) or average prediction (regression) of the individual trees. Let T denote the number of trees in the forest. For classification, the random forest combines the predictions $\hat{y}_{(t)}$ of each tree t by voting: $\hat{y} = \text{mode}(\hat{y}_{(1)}, \dots, \hat{y}_{(t)})$.
- **XGBoost (Extreme Gradient Boosting).** XGBoost is an optimized gradient-boosting library known for its speed and performance. It sequentially builds trees, where each subsequent tree corrects errors made by the previous one. XGBoost minimizes the loss function $L(\phi)$ iteratively by adding weak learners h_t to the model: $\hat{y} = \phi(x) = \sum_{t=1}^T h_t(x)$, where T is the number of boosting rounds.
- **AdaBoost.** AdaBoost is another ensemble learning method that combines multiple weak classifiers to create a strong classifier. It adjusts the weights of incorrectly classified instances so that subsequent classifiers focus more on difficult cases. At each iteration t , AdaBoost updates the weights $D_t(i)$ of training instances and computes the model weight α_t based on the classification error. The final classifier is a weighted sum: $\hat{y} = \text{sign}(\sum_{t=1}^T \alpha_t h_t(x))$.
- **GBDT.** GBDT builds trees sequentially, where each tree attempts to correct errors made by the previous one. Unlike AdaBoost, which adjusts instance weights, GBDT fits each new tree to the residual errors of the current model predictions. GBDT minimizes the loss function by adding new trees that approximate the negative gradient of the loss function for the ensemble model: $\hat{y} = \sum_{t=1}^T h_t(x)$, where h_t is the t -th decision tree.
- **HGBDT.** HGBDT is an optimized version of GBDT that uses histograms to discretize continuous features, reducing memory usage and speeding up training. Similar to GBDT, HGBDT constructs an ensemble model by iteratively adding decision trees that minimize the loss function: $\hat{y} = \sum_{t=1}^T h_t(x)$, where each tree h_t is trained on histogram-based feature representations.

3.3. Local Explanation for SQL Injection Model Decision Based on SHAP and LIME

In the context of SQL injection detection, understanding why a machine learning model classifies an SQL query as normal or malicious is crucial to ensuring the reliability and robustness of the system. Local explanation techniques, such as SHAP and LIME, are widely used to explain individual predictions. These methods help elucidate the model decision-making process by providing feature attributions for a specific instance. This section details how SHAP and LIME can be used to explain the behavior of machine learning models to detect SQL injection attacks.

3.3.1. SHAP: SHapley Additive exPlanations

SHAP is a method based on cooperative game theory, specifically Shapley values, which provides a theoretically grounded approach to explain the prediction of a model. Shapley values calculate the marginal contribution of each feature to the prediction of the model by considering all possible combination of features.

Mathematically, the Shapley value for a feature i is defined as:

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

where, N is the set of all features, $S \subseteq N \setminus \{i\}$ is a subset of features excluding feature i , $f(S)$ is the model's prediction when only the features in S are considered, $f(S \cup \{i\})$ is the model's prediction when feature i is added to S , ϕ_i represents the contribution of feature i to the model's prediction.

For our SQL injection detection model, SHAP values can be computed for each feature in the SQL query, such as SQL keywords (e.g., 'SELECT', 'UNION'), query length, and query structure. The SHAP values indicate how each feature contributes to classifying the query as either normal or malicious.

To visualize these attributions, we use SHAP's `force_plot`, which provides an intuitive representation of the forces that push the model's prediction towards a certain class. The force plot for two cases—one for a normal query and one for an SQL injection attack—illustrates how specific features (e.g. keywords or anomalous patterns) drive the decision.

3.3.2. LIME: Local Interpretable Model-Agnostic Explanations

LIME provides local explanations by approximating the decision boundary of the original model with an interpretable surrogate model, such as a linear regression or a decision tree, in the neighborhood of the instance to be explained. LIME works by generating perturbed samples around the instance of interest and fitting a local model to explain how the original model predictions change in response to those perturbations.

The prediction for an instance x is explained by solving the following minimization problem:

$$\zeta(x) = \arg \min_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g) \quad (2)$$

where, f is the original model, g is the interpretable surrogate model from the class of models G (e.g., a linear model), $\mathcal{L}(f, g, \pi_x)$ is the loss function that measures how well the surrogate model g approximates the original model f in the local neighborhood defined by π_x , $\Omega(g)$ is a complexity term that ensures the surrogate model remains interpretable.

For SQL injection detection, we apply `LimeTabularExplainer` to both normal and malicious SQL queries. `LimeTabularExplainer` generates new data points by randomly perturbing the original query and then examines how the model responds to these variations. Build a local interpretable model to approximate the decision-making process for that specific query.

4. Experiment

Our models are developed in Python, utilizing widely used machine learning libraries such as Scikit-learn. The experiments are performed on a computer with the following specifications: Intel Core i7-10700K CPU @ 3.80GHz processor, 64GB RAM, and Windows 10 operating system.

4.1. Experiment Setting

Table 1 summarizes the key hyperparameters used across multiple models, including `XGBClassifier`, `RandomForestClassifier`, `DecisionTreeClassifier`, `GradientBoostingClassifier`, `HistGradientBoostingClassifier`, and `AdaBoostClassifier`. The table highlights essential hyperparameters such as the number of estimators, learning rate, loss function, and max depth, which significantly impact model training and performance. These configurations provide a clear overview of the settings employed for each classifier in the experiments.

Table 1. Hyperparameters for ensemble learning and boosting models

Model Name	Hyperparameter	Value
XGBClassifier	enable_categorical	False
	n_estimators	100
RandomForestClassifier	n_estimators	100
	criterion	gini
	max_depth	None
	min_samples_split	2
	min_samples_leaf	1
	max_features	auto
DecisionTreeClassifier	criterion	gini
	splitter	best
	max_depth	None
	min_samples_split	2
	random_state	42
GradientBoostingClassifier (GBDT)	loss	deviance
	learning_rate	0.1
	n_estimators	100
	subsample	1.0
	criterion	friedman_mse
	min_samples_split	2
	min_samples_leaf	1
	max_depth	3
	max_features	None
random_state	42	
HistGradientBoostingClassifier (HGBDT)	loss	auto
	learning_rate	0.1
	max_iteractions	100
	max_leaf_nodes	31
	max_depth	None
	min_samples_leaf	20
	l2_regularization	0.0
	max_bins	255
	early_stopping	auto
	scoring	loss
random_state	42	
AdaBoostClassifier	n_estimators	50
	learning_rate	1.0
	algorithm	SAMME.R
	base_estimator	DecisionTreeClassifier
	random_state	None
	estimator_params	()

4.2. Evaluation Metrics

To assess the performance of our binary classifier, we use several evaluation metrics, including the confusion matrix, precision, recall, F1 score, accuracy and the Receiver Operating Characteristic (ROC) curve. Each of these metrics provides valuable insights into the model's performance and its ability to classify instances correctly.

The confusion matrix is a table that summarizes the performance of a classification algorithm. It displays the counts of true positive (TP), false positive (FP), true negative (TN), and false negative (FN) predictions as follows:

The precision known as the positive predictive value measures the precision of the positive predictions. It is calculated as:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (3)$$

Recall, also known as sensitivity or true positive rate, measures the ability of the model to identify positive instances. It is defined as:

$$\text{Recall} = \frac{TP}{TP + FN} \quad (4)$$

The F1 score is the harmonic mean of precision and recall, providing a single metric that balances both concerns. It is computed as:

$$\text{F1 Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (5)$$

Accuracy measures the overall correctness of the model by calculating the ratio of correctly predicted instances to the total instances:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (6)$$

Receiver Operating Characteristic (ROC) Curve: The ROC curve plots the True Positive Rate (TPR) against the False Positive Rate (FPR) at various threshold settings. The TPR and FPR are defined as:

$$\text{TPR} = \frac{TP}{TP + FN} \quad (7)$$

$$\text{FPR} = \frac{FP}{FP + TN} \quad (8)$$

The area under the ROC curve (AUC-ROC) is a single scalar value that summarizes the performance of the classifier across all thresholds, where a higher AUC indicates better performance.

These evaluation metrics collectively provide a comprehensive view of the performance of the model, allowing us to make informed decisions regarding its effectiveness in the classification task at hand.

4.3. Confusion Matrix Results

This section presents the confusion matrix results for our models in classifying SQL injection attacks and non-attacks. The Decision Tree, Random Forest, and AdaBoost models exhibit the best detection performance, closely followed by the GBDT model. In contrast, XGBoost demonstrates the least accuracy, misclassified 1,424 non-attack samples as an attack (Figure 3(c)).

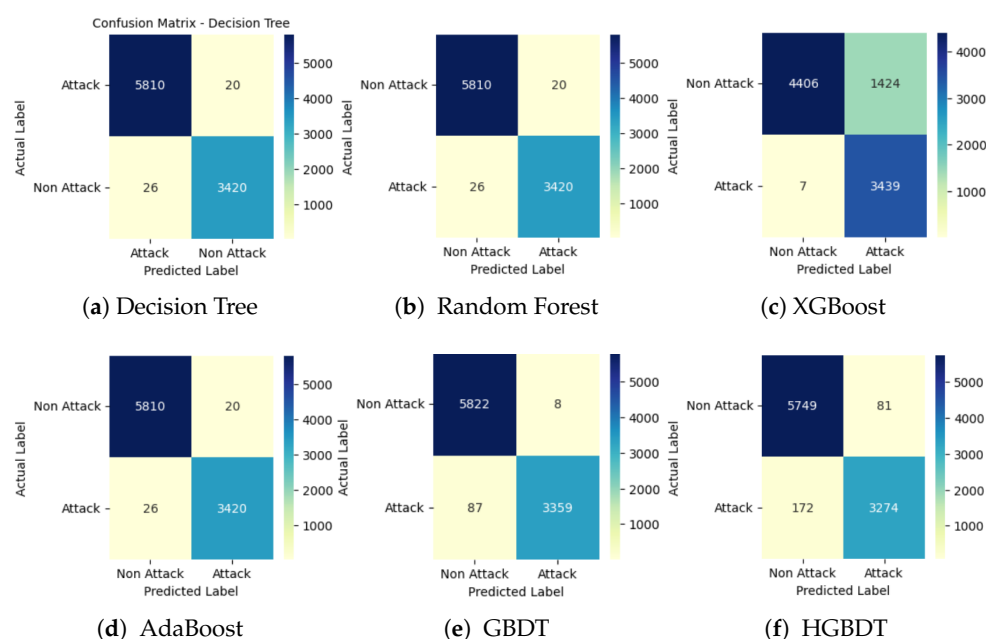


Figure 3. Confusion matrix results of variant decision tree models.

4.4. ROC Curve Results

We analyze the diagnostic performance of a binary classifier system by varying its discrimination threshold. The ROC curve plots the true positive rate (TPR) against the false positive rate (FPR) across various threshold settings. The area under the ROC curve (AUC-ROC) serves as a concise metric of the classifier's effectiveness, where a higher AUC-ROC signifies the superior discriminatory capability of the model.

Figure 4 presents the ROC measurements for six models. The results indicate that XGBoost achieves an AUC-ROC of 88%, while the remaining models achieve robust AUC-ROC values exceeding 97%.

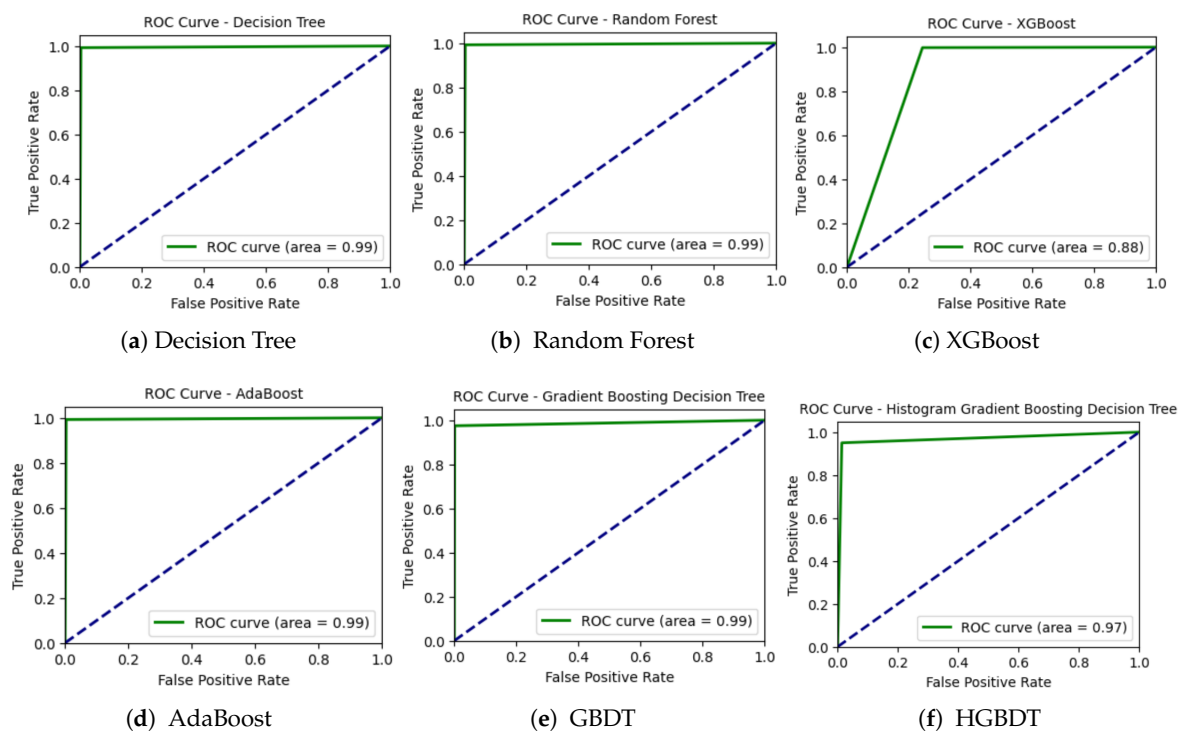


Figure 4. ROC curve results of variant decision tree models.

4.5. Other Performance Evaluation Matrix Results

This section comprehensively analyzes the performance of our models using several evaluation metrics, including precision, recall, accuracy, and the F1 score. The evaluation results, summarized in Table 2, illustrate the effectiveness of each model in detecting SQL injection attacks.

Table 2. Detailed classification evaluation results.

Model	Precision	Recall	Accuracy	F1	ROC
Decision Tree	0.9942	0.9925	0.9950	0.9933	0.99
Random Forest	0.9942	0.9925	0.9950	0.9933	0.99
XGBoost	0.7072	0.9980	0.8457	0.8278	0.88
AdaBoost	0.9942	0.9925	0.9950	0.9933	0.99
GBDT	0.9976	0.9748	0.9898	0.9861	0.99
HGBDT	0.9759	0.9501	0.9727	0.9628	0.97

Decision Tree, Random Forest, and AdaBoost models demonstrate exceptional performance, each achieving a precision of 99.42%, recall of 99.25%, accuracy of 99.50%, and an F1 score of 99.33%. These results indicate that these models are highly effective in identifying SQL injection attacks while maintaining a low rate of false positives.

In contrast, XGBoost, while achieving a high recall of 99.80%, has a lower precision of 70.72%, resulting in an overall accuracy of 84.57% and an F1 score of 82.78%. This suggests that although XGBoost is good at detecting almost all attack instances, it produces more false positives than the other models.

The GBDT model shows strong performance with a precision of 0.9976, recall of 97.48%, accuracy of 98.98%, and an F1 score of 98.61%. This indicates that GBDT is also a reliable model, though slightly less accurate than Decision Tree, Random Forest, and AdaBoost.

Finally, the HGBDT model, while still effective, lags slightly behind the top performers with a precision of 97.59%, a recall of 95.01%, a precision of 97.27%, and an F1 score of 96.28%. These metrics show that HGBDT is a solid model but not as robust as the other ensemble methods evaluated.

Our analysis highlights that ensemble methods like Random Forest and AdaBoost, along with the Decision Tree model, provide the most reliable performance detecting SQL injection attacks. These findings emphasize the importance of model selection and tuning in the development of effective cybersecurity solutions.

4.6. Comparison with Existing Methods

Table 3 presents a performance comparison between our approach and several existing methods for SQL injection detection.

Table 3. Comparison of classification performance with other prior methods.

Model	Precision	Recall	Accuracy	F1	ROC
Logistic Regression [32]	-	-	99.3	-	-
SQLIA [33]	97.4	99.7	98	98.5	99.9
CNN [34]	-	-	99.6	-	-
Naive Bayes [35]	94.19	-	98.33	97.00	97.71
Ours	99.42	99.25	99.50	99.33	99

Our model demonstrates competitive or superior performance across multiple metrics. With an accuracy of 99.50%, our method slightly outperforms the CNN model (99.6%) [34] and Logistic Regression (99.3%) [32]. The F1 score of our model (99.33%) surpasses that of SQLIA (98.5%) [33] and Naive Bayes (97.00%) [35]. While SQLIA reports a marginally higher ROC (99.9%) compared to our 99%, our model maintains a better overall balance across all metrics. These results underscore the effectiveness of our approach in the context of current SQL injection detection techniques.

4.7. Local Explanation Results for Model Decision

We explain the decisions made by our ensemble and the boost models using SHAP and LIME values on two types of data samples: attack samples and normal samples. For SHAP, we use a force graph to display the SHAP values. For LIME, we visualize prediction probabilities and explain the values that lead to each model's decision regarding data detection.

First, we visualize the interpretation of SHAP values for ensemble models, including Decision Tree and Random Forest, as shown in Figure 5. Second, we visualize the interpretation of SHAP values for boosting models, including XGBoost, AdaBoost, GBDT, and HGBDT, as illustrated in Figure 6.

Based on the SHAP visualized results, we observe that a positive SHAP value (represented in red) indicates that the model's decision is an attack classification, whereas a negative SHAP value (represented in blue) signifies a normal classification.

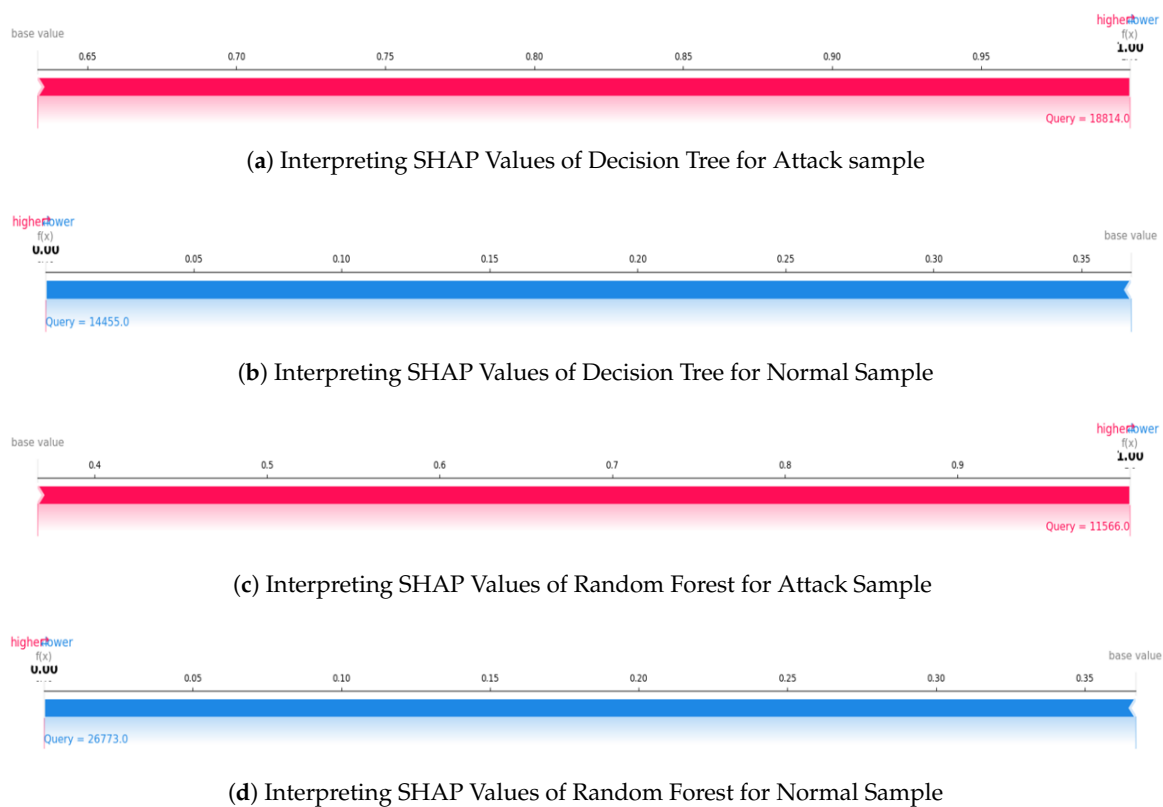
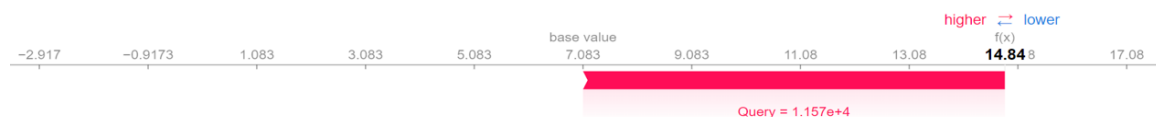
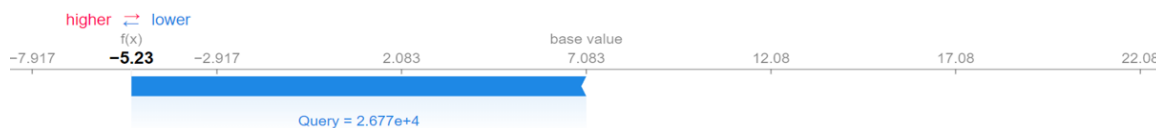


Figure 5. Interpreting SHAP Values in Ensemble Models for Normal and Attack Samples

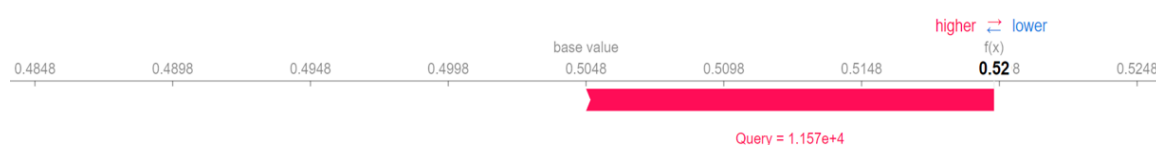
Next, we visualize the interpretation of LIME values for six models, as shown in Figures 7 to 12. Based on the visualized results from LIME, we recognize that a positive LIME value (represented in green) indicates that the model's decision is an attack classification, while a negative LIME value (represented in red) signifies a normal classification.



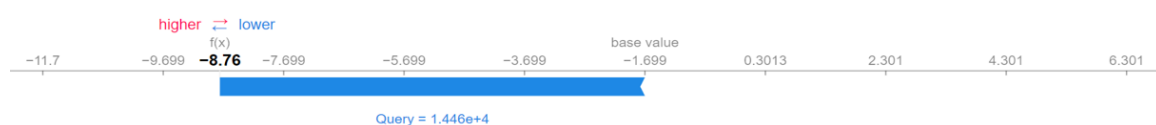
(a) Interpreting SHAP Values of XGBoost for Attack Sample



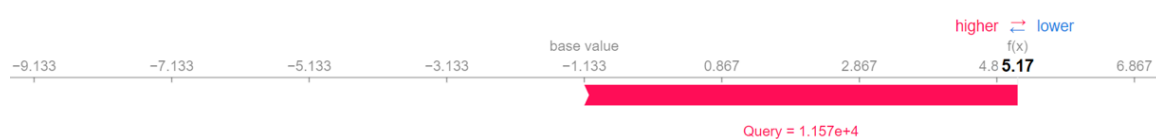
(b) Interpreting SHAP Values of XGBoost for Normal Sample



(c) Interpreting SHAP Values of AdaBoost for Attack Sample



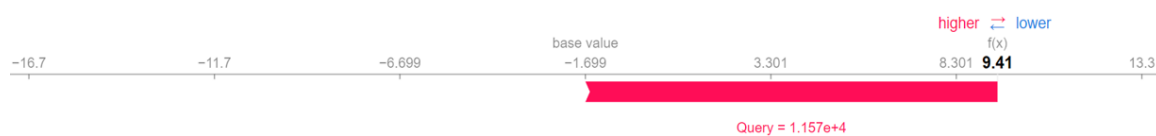
(d) Interpreting SHAP Values of AdaBoost for Normal Sample



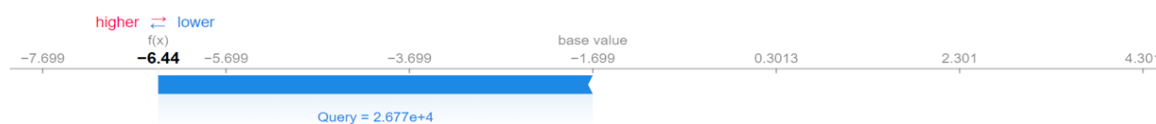
(e) Interpreting SHAP Values of GBDT for Attack Sample



(f) Interpreting SHAP Values of GBDT for Normal Sample

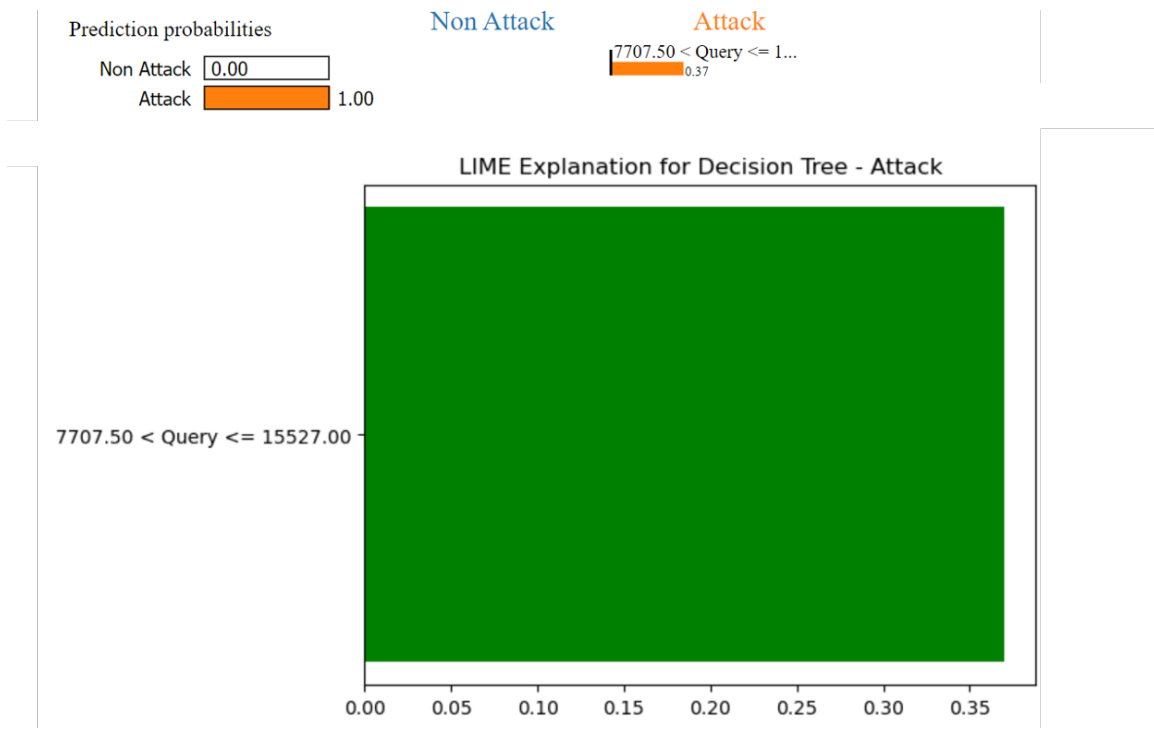


(g) Interpreting SHAP Values of HGBDT for Attack Sample

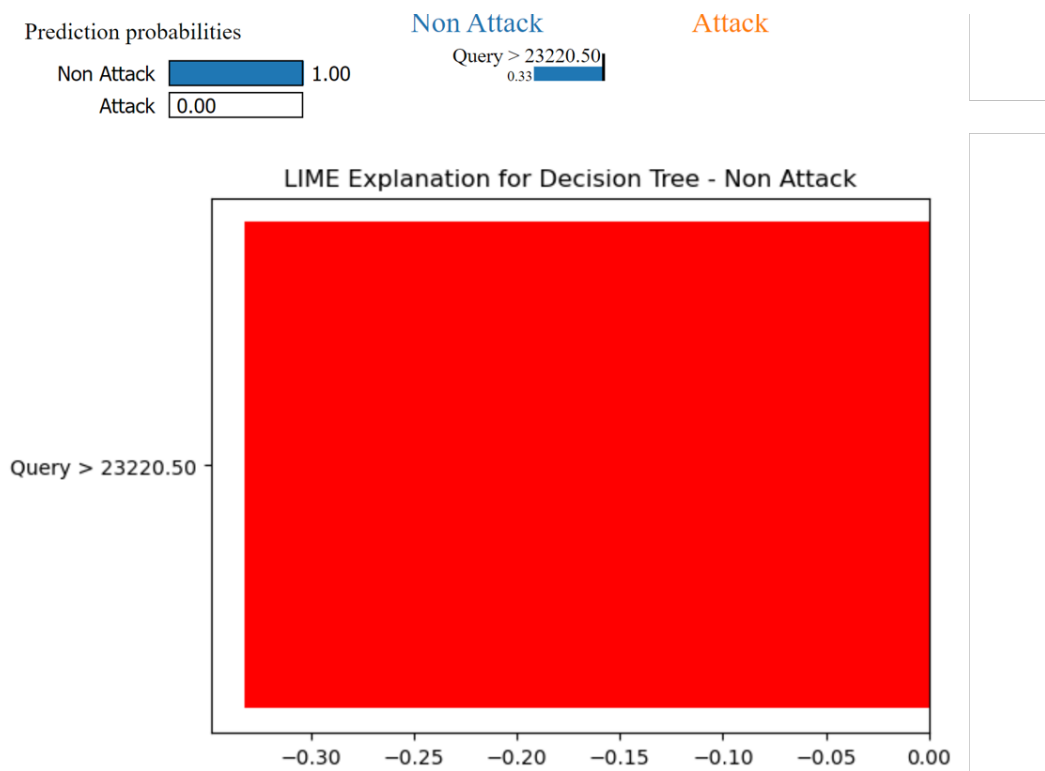


(h) Interpreting SHAP Values of HGBDT for Normal Sample

Figure 6. Interpreting SHAP Values in Boosting Models for Normal and Attack Samples

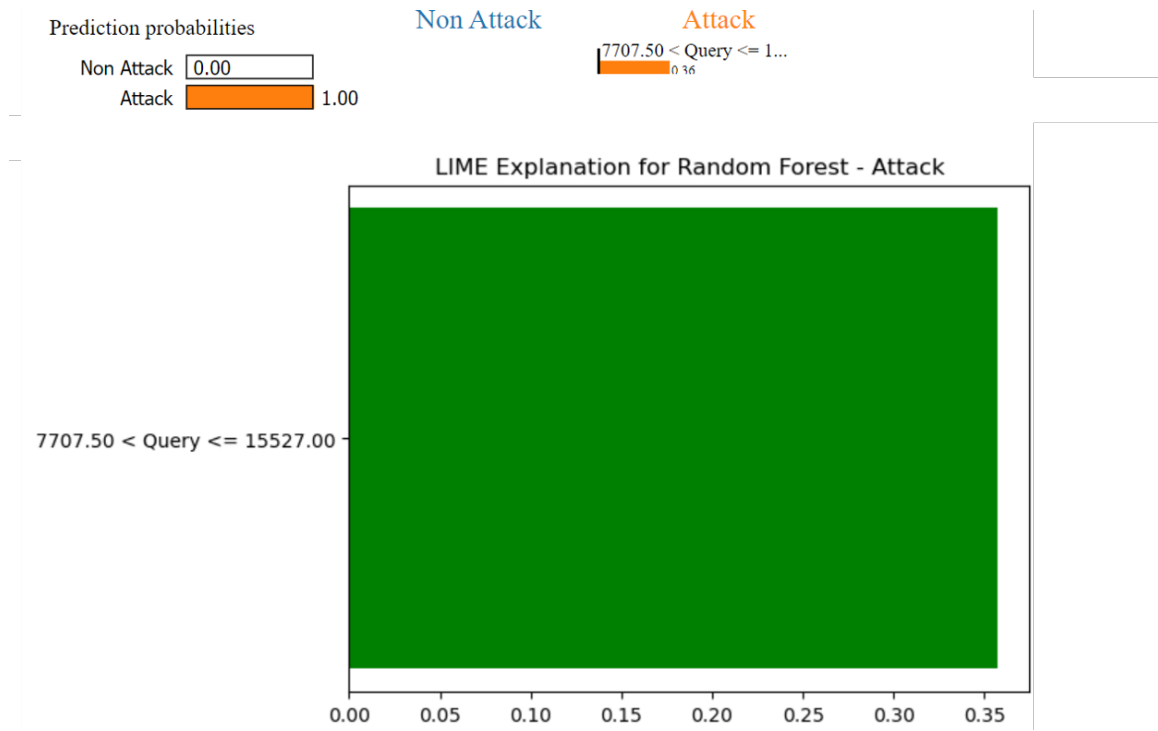


(a) Interpreting LIME Values of Decision Tree for Attack Sample

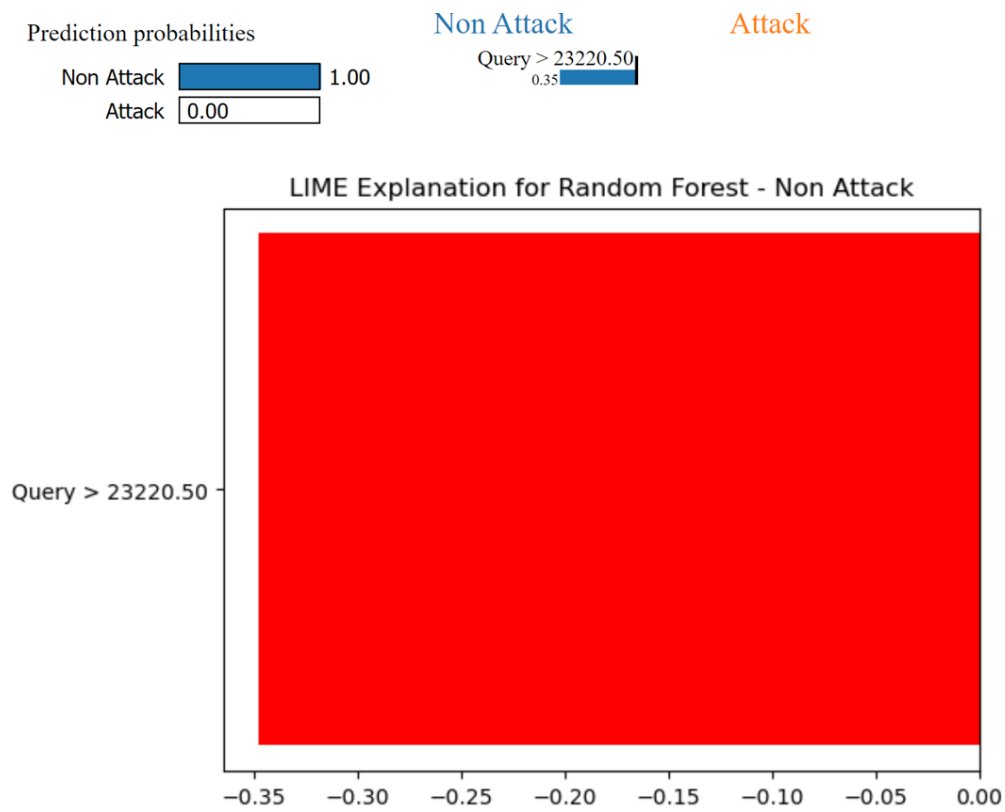


(b) Interpreting LIME Values of Decision Tree for Normal Sample

Figure 7. Interpreting SHAP Values in Decision Tree Model for Normal and Attack Samples

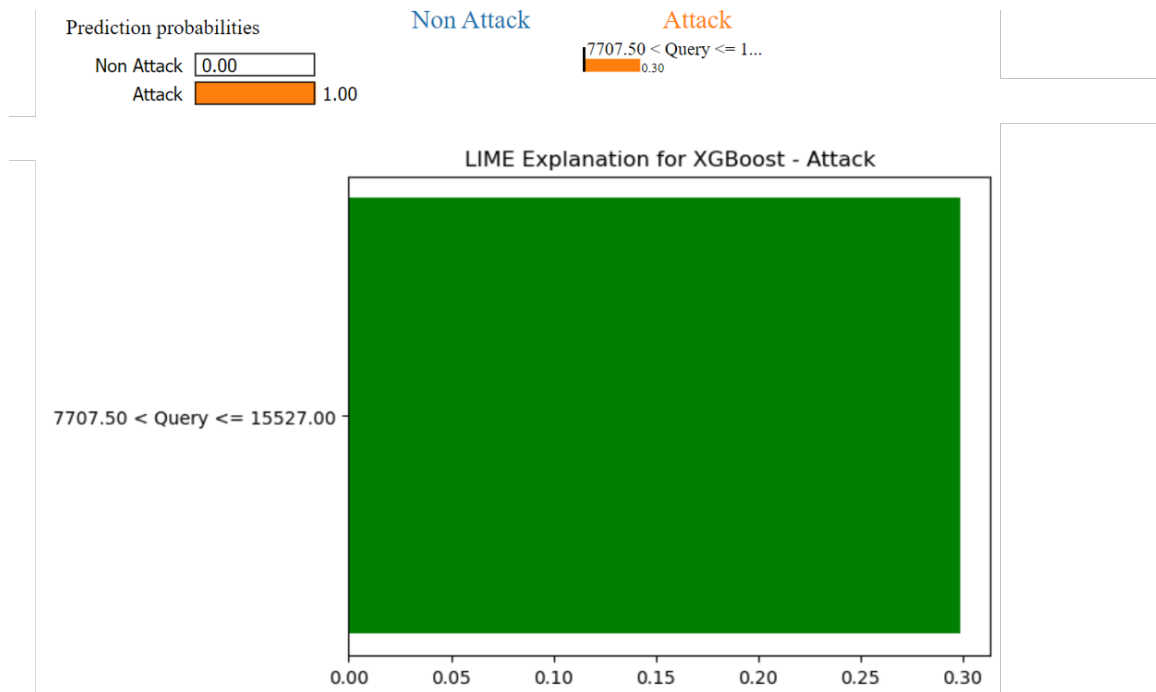


(a) Interpreting LIME Values of Decision Tree for Attack Sample

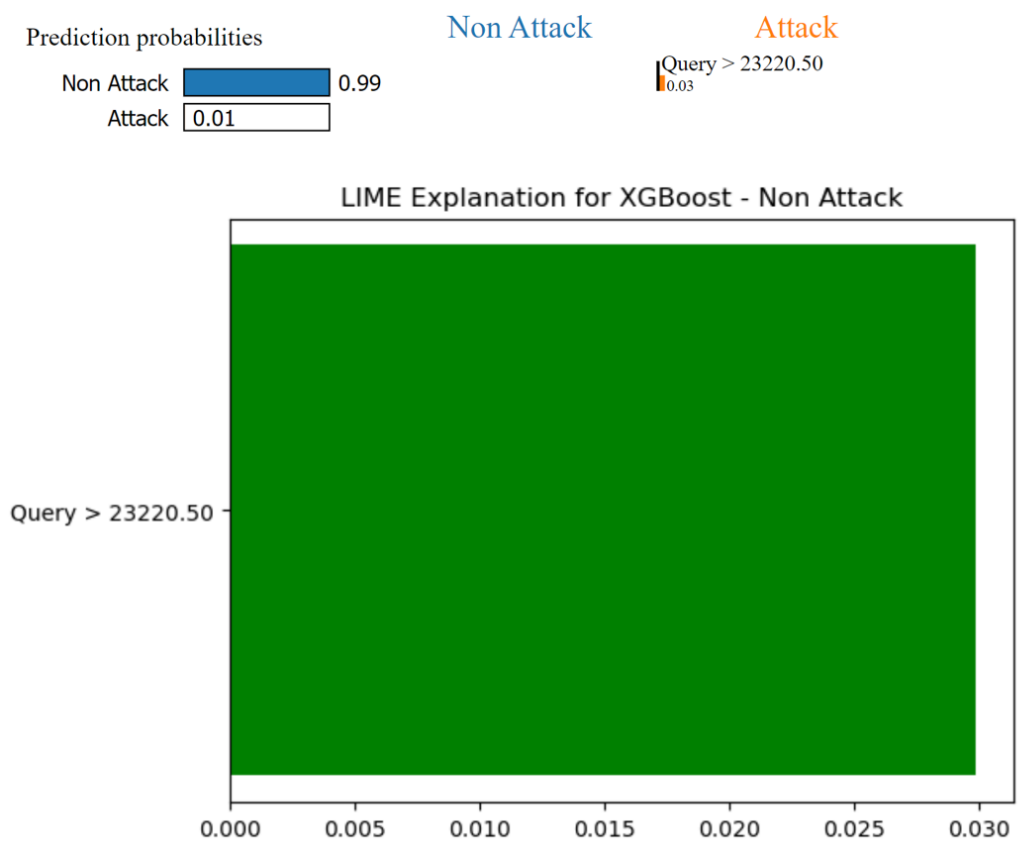


(b) Interpreting LIME Values of Decision Tree for Normal Sample

Figure 8. Interpreting SHAP Values in Random Forest Model for Normal and Attack Samples

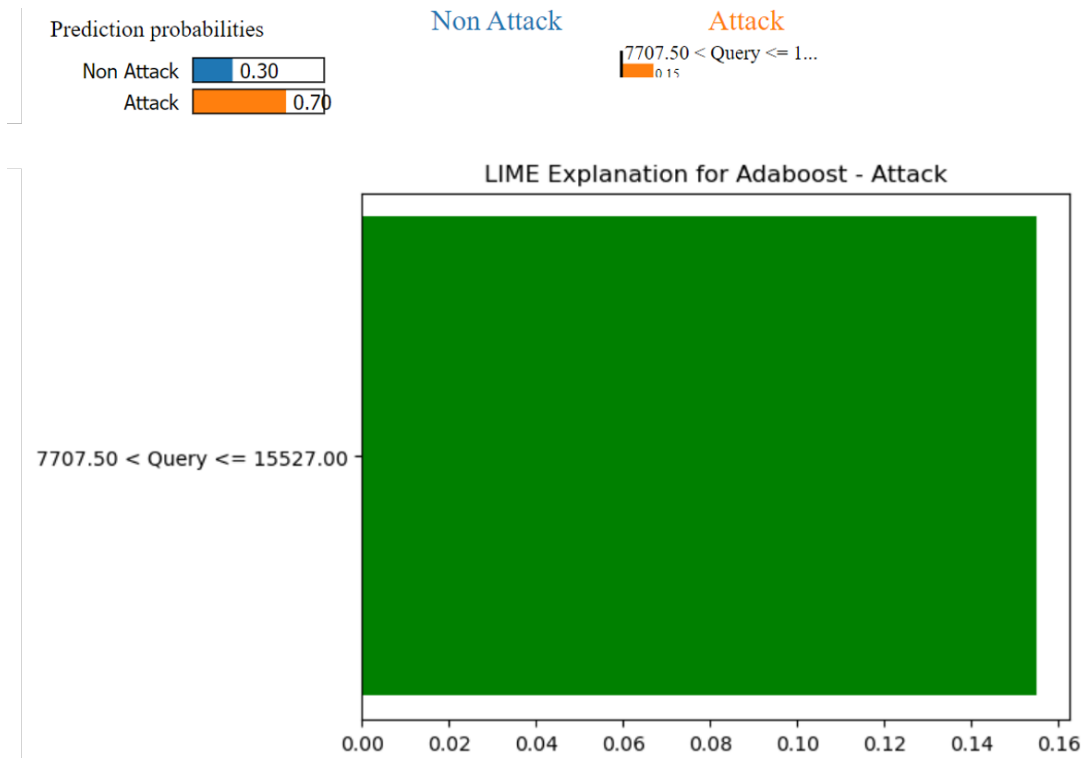


(a) Interpreting LIME Values of XGBoost for Attack Sample

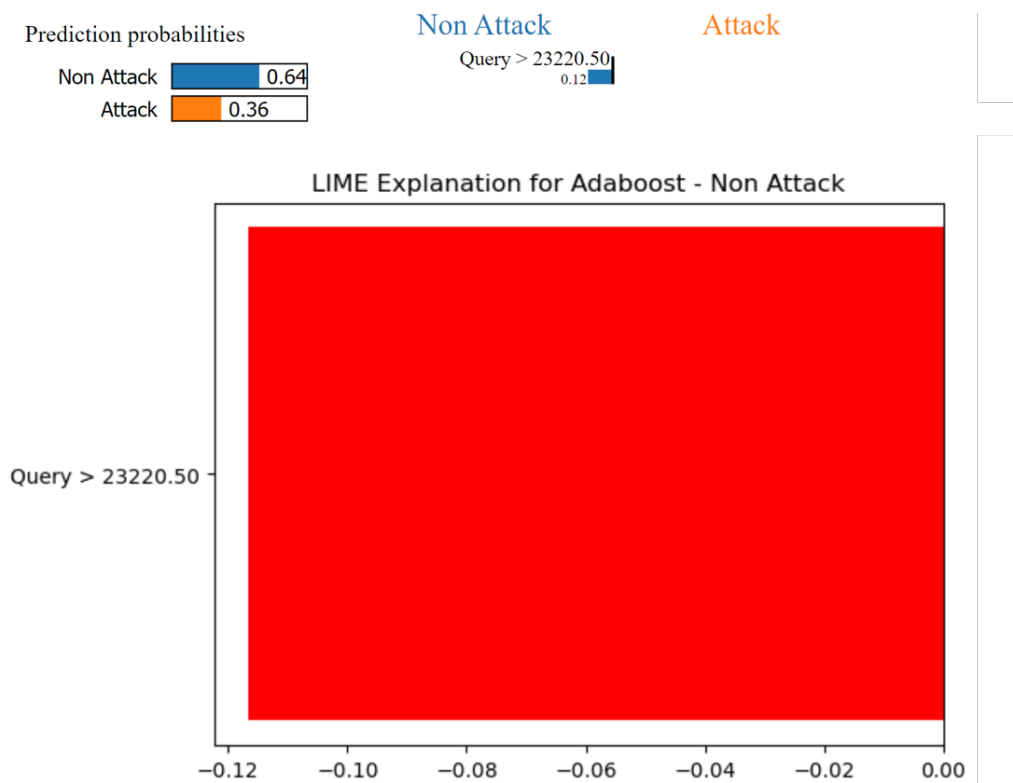


(b) Interpreting LIME Values of XGBoost for Normal Sample

Figure 9. Interpreting SHAP Values in XGBoost Model for Normal and Attack Samples

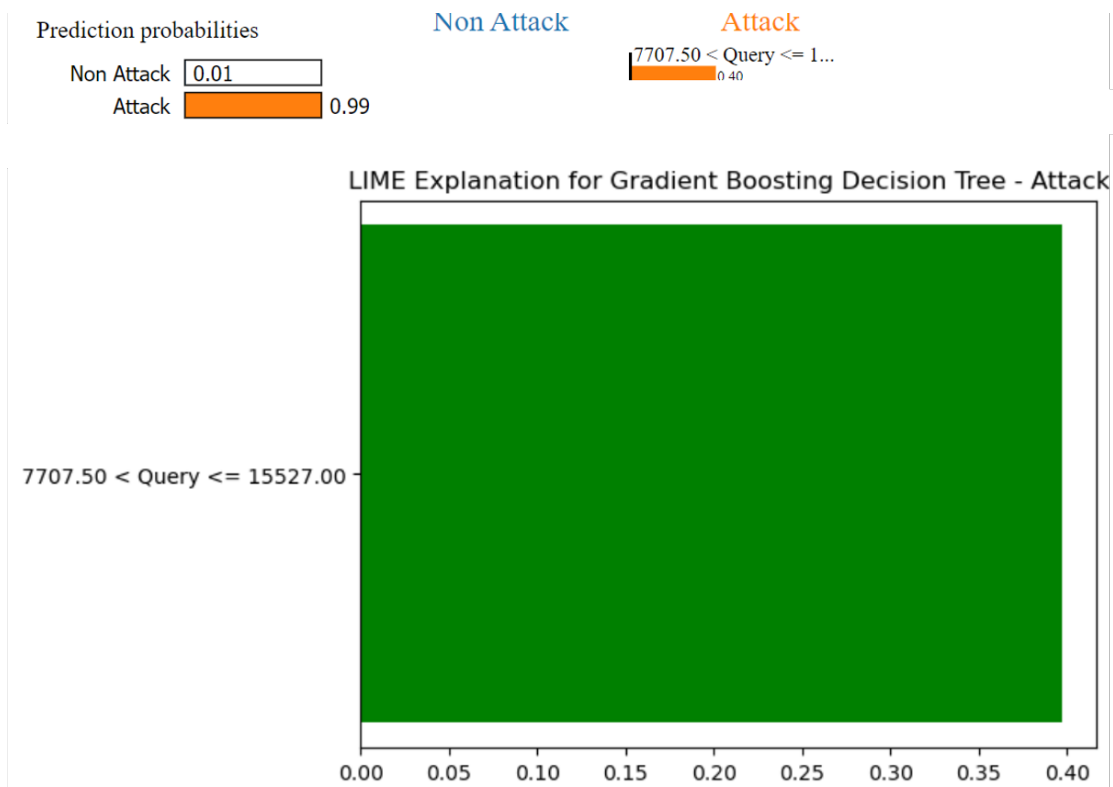


(a) Interpreting LIME Values of AdaBoost for Attack Sample

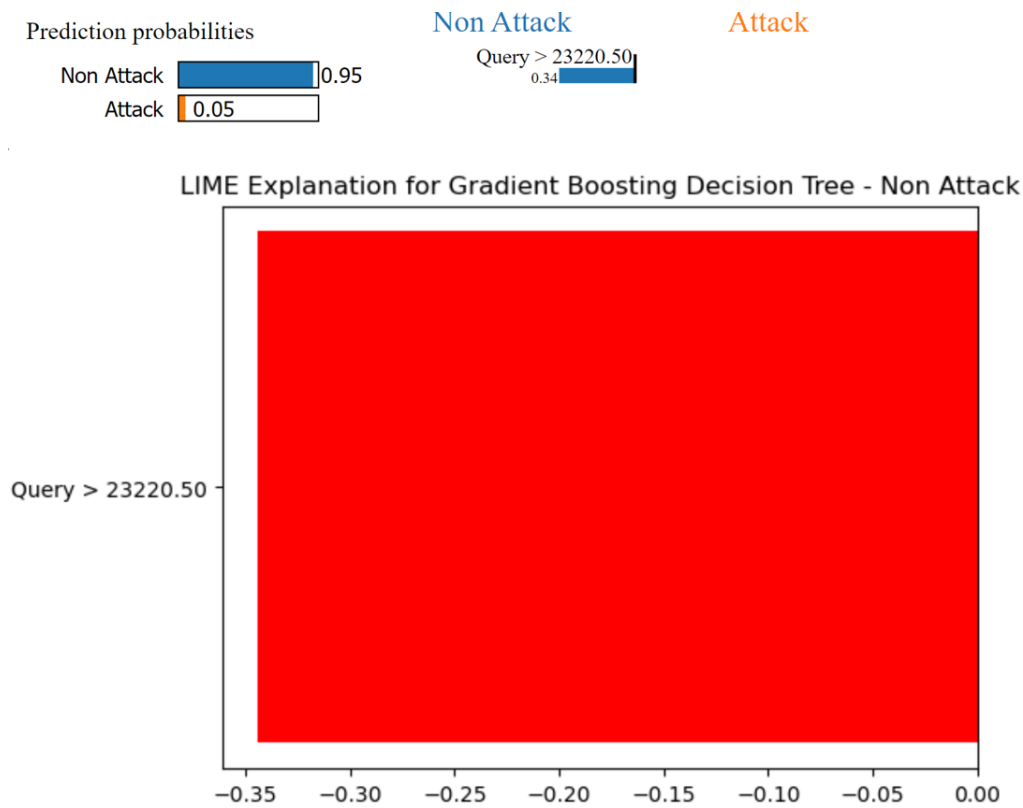


(b) Interpreting LIME Values of AdaBoost for Normal Sample

Figure 10. Interpreting SHAP Values in AdaBoost Model for Normal and Attack Samples

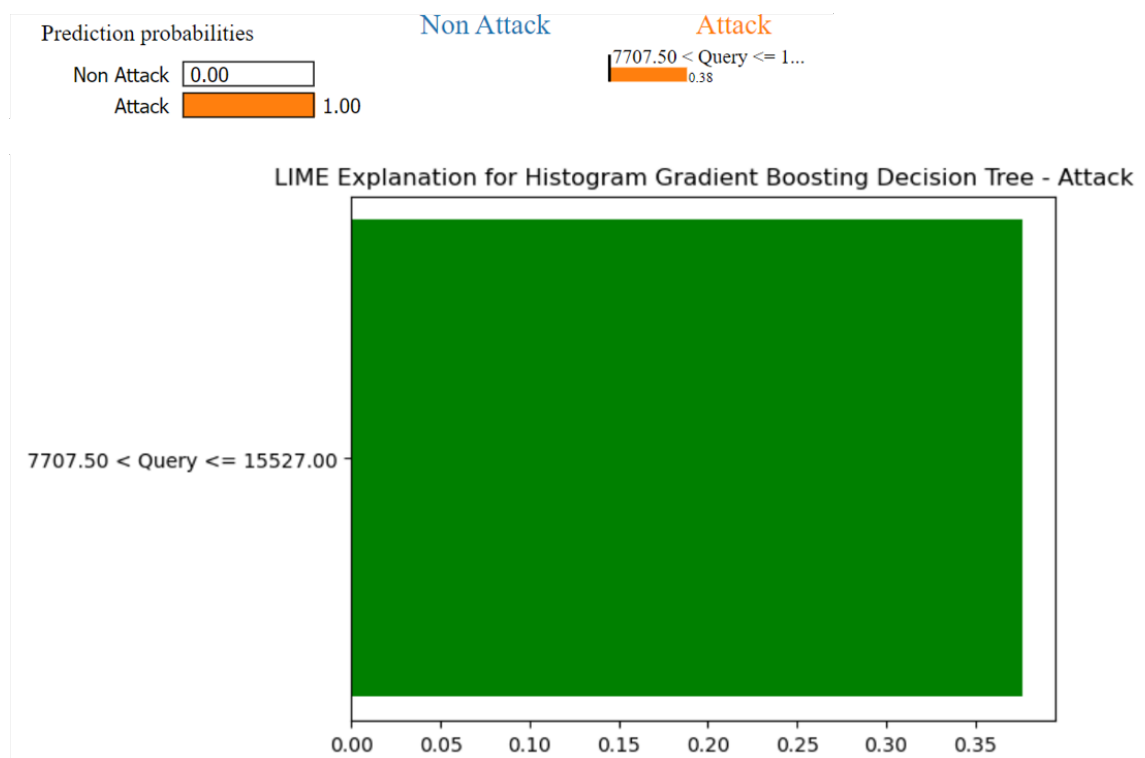


(a) Interpreting LIME Values of GBDT for Attack Sample

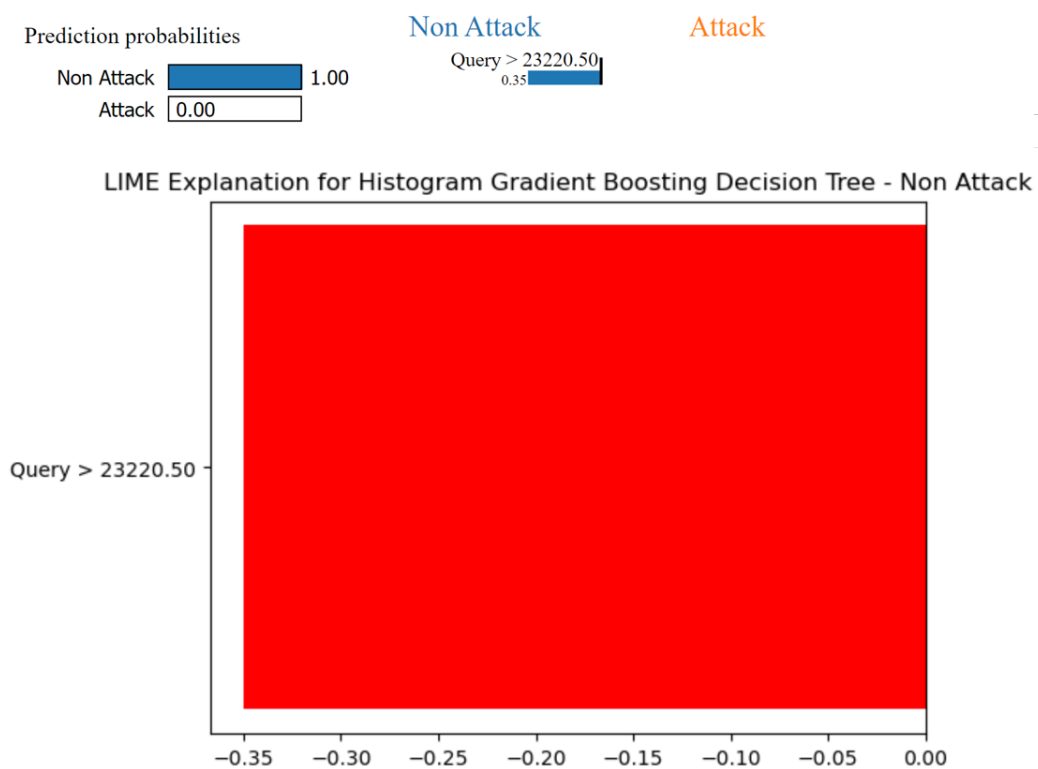


(b) Interpreting LIME Values of GBDT for Normal Sample

Figure 11. Interpreting SHAP Values in GBDT Model for Normal and Attack Samples



(a) Interpreting LIME Values of HGBDT for Attack Sample



(b) Interpreting LIME Values of HGBDT for Normal Sample

Figure 12. Interpreting SHAP Values in HGBDT Model for Normal and Attack Samples

5. Discussion

This study's experimental findings reveal the superior performance of Decision Tree, Random Forest, and AdaBoost models in detecting SQL injection attacks, with outstanding metrics including

99.50% accuracy, 99.33% F1 score, and 99% AUC-ROC. These results not only demonstrate the models' reliability but also offer insights into the nature of SQL injection detection.

The success of tree-based and ensemble methods in this context can be attributed to several factors. Firstly, SQL injection attacks often involve specific patterns or structures in the input data. Tree-based models excel at capturing these hierarchical relationships, allowing them to distinguish between benign and malicious SQL queries effectively. Secondly, the ensemble methods (Random Forest and AdaBoost) leverage the power of multiple models, which is particularly beneficial in cybersecurity applications where attack vectors can be diverse and evolving.

Comparing our results to existing literature, our top-performing models achieve slightly higher accuracy than some previous studies. For instance, [35] reported 97.71% F1 score using Naive Bayes, while our best models exceed 99% F1 score. This improvement, although marginal, could translate to significant real-world benefits in reducing false positives and negatives in SQL injection detection systems.

The high performance of simpler models like Decision Trees is particularly noteworthy. While complex deep learning approaches have gained popularity in various domains, our results suggest that simpler, interpretable models can equally effectively detect SQL injection. This finding has important implications for practical implementations, as these models are computationally less intensive and more accessible to deploy and maintain in real-time detection systems.

5.1. Insights from Local Explanation Results

The application of SHAP and LIME for local interpretability provides additional valuable insights into the decision-making process of our models. The SHAP force plots and LIME explanations clearly identify the features that contribute the most to each model's classification of attack or normal samples.

From the SHAP visualizations, it is evident that certain features have a strong positive impact on predicting SQL injection attacks, while others consistently support normal behavior classifications. These findings highlight that specific input patterns, such as certain SQL query components, are consistently recognized as attack indicators by models like Decision Tree and Random Forest. This consistency across various samples demonstrates the reliability of our models' pattern-recognition capabilities.

LIME further reinforces these observations by showing the prediction probabilities and explaining feature contributions in a clear and interpretable manner. The alignment between the SHAP and LIME results strengthens the confidence in the robustness of the models and their transparency, making them more trustworthy for practical cybersecurity applications.

The insights from these local explanation techniques also suggest potential areas for model improvement. For example, in cases where certain features contributed inconsistently across different boosting models (such as GBDT and HGBDT), this may indicate opportunities for better feature engineering or adjustments in model training strategies to enhance performance.

5.2. Limitations and Future Work

However, despite these promising results, it is crucial to address some limitations of this study. The slightly lower performance of gradient boosting methods (GBDT, HGBDT, XGBoost) compared to tree-based and ensemble models warrants further investigation. The explanation results for these models indicated less consistent feature contributions, suggesting that these models might require more refined hyperparameter tuning or feature selection tailored to the characteristics of SQL injection data.

Moreover, while our models performed well in detecting known attack patterns, future work should explore their robustness in real-world, adversarial scenarios where attackers actively attempt to evade detection. Techniques such as adversarial training and further model fine-tuning could help improve the models' resilience in these more challenging environments.

In conclusion, the integration of local explanation methods such as SHAP and LIME not only enhances the transparency and trustworthiness of our models but also provides actionable insights for further improving SQL injection detection systems. These findings lay the groundwork for future research on refining feature engineering, enhancing model robustness, and testing against evolving attack vectors.

6. Conclusion

This paper presents a comparative analysis of various decision models for detecting SQL injection attacks, with a focus on Decision Tree, Random Forest, XGBoost, AdaBoost, GBDT, and HGBDT. Among these, the Decision Tree, Random Forest, and AdaBoost models exhibit superior performance, achieving high accuracy in identifying SQL injection attacks. Their robust performance, with accuracy metrics exceeding 99%, underscores their reliability and suitability for real-world deployment in enhancing cybersecurity defenses.

While gradient-boosting methods such as GBDT and HGBDT demonstrate competitive results, they do not outperform the simpler ensemble methods like Random Forest and AdaBoost. Notably, XGBoost, despite its popularity, shows limitations in this specific application, suggesting that additional refinement or alternative techniques may be required to optimize its performance for SQL injection detection.

Incorporating local explanation methods such as SHAP and LIME has further enriched our understanding of model behavior, providing transparency into the decision-making processes of our models. These local explanations reveal how certain features consistently contribute to attack and normal classifications, thereby increasing trust in the models' predictions. This interpretability is crucial in cybersecurity contexts, where understanding the reasoning behind attack detection is as important as the accuracy of the detection itself.

Looking ahead, future research will leverage the insights from these models to develop a proactive tool for generating SQL injection attacks using text generation techniques. This tool aims to provide deeper insight into the nature of SQL vulnerabilities, helping cybersecurity professionals preemptively address potential threats. In addition, integrating our top-performing models with real-time monitoring systems will allow for the assessment of their performance in dynamic, adversarial environments where attackers continually adapt their methods to evade detection.

Finally, exploring hybrid approaches that combine the strengths of various classifiers could further enhance detection rates and resilience against sophisticated SQL injection techniques. As SQL injection attacks continue to evolve, adaptive and interpretable models, like those presented in this study, will be essential for maintaining robust security in database-driven applications.

Author Contributions: Conceptualization, L.T.T.H. and H.Y.; methodology, L.T.T.H.; software, L.T.T.H. and H.Y.; validation, L.T.T.H. and K.H.; formal analysis, H.Y, C. C.; investigation, H.Y and C.C; resources, C.C, W.W.R, and P.C.S.D; data curation, H.Y, and C.C; writing—original draft preparation, L.T.T.H, W.W.R, and P.C.S.D; writing—review and editing, L.T.T.H. and K.H.; visualization, L.T.T.H, H.Y, C.C, W.W.R and P.C.S.D; supervision, L.T.T.H and K.H; project administration, K.H; funding acquisition, K.H. All authors have read and agreed to the published version of the manuscript.

Funding: This research was supported by the MSIT(Ministry of Science and ICT), Korea, under the ITRC(Information Technology Research Center) support program(IITP-2023-2020-0-01797) supervised by the IITP(Institute for Information & Communications Technology Planning & Evaluation).

Institutional Review Board Statement: Not applicable

Informed Consent Statement: Not applicable

Data Availability Statement: Not applicable

Conflicts of Interest: The authors declare no conflicts of interest.

References

1. OWASP Top Ten. The Ten Most Critical Web Application Security Risks. Retrieved from <https://owasp.org/www-project-top-ten/>, last accessed by 2024/05/02.
2. Halfond, W. G., Viegas, J., & Orso, A.: A classification of SQL-injection attacks and countermeasures. In Proceedings of the IEEE international symposium on secure software engineering (Vol. 1, pp. 13-15). Piscataway, NJ: IEEE, . (2006, March).
3. Jacob, I., & Pirnau, M. (2020). SQL INJECTION ATTACKS AND VULNERABILITIES. *Journal of Information Systems & Operations Management*, 68-81.
4. Su, Z., & Wassermann, G. The essence of command injection attacks in web applications. *Acm Sigplan Notices*, 41(1), 372-382, (2006).
5. Demilie, W. B., & Deriba, F. G. (2022). Detection and prevention of SQLI attacks and developing compressive framework using machine learning and hybrid techniques. *Journal of Big Data*, 9(1), 124.
6. Shar, L. K., & Tan, H. B. K. (2013). Predicting SQL injection and cross site scripting vulnerabilities through mining input sanitization patterns. *Information and Software Technology*, 55(10), 1767-1780.
7. Nair, S. S. (2024). Securing Against Advanced Cyber Threats: A Comprehensive Guide to Phishing, XSS, and SQL Injection Defense. *Journal of Computer Science and Technology Studies*, 6(1), 76-93.
8. Anley, C.: Advanced SQL injection in SQL server applications, NGS Software Insight Security Research (NISR) (2002).
9. Chandola, V., Banerjee, A., & Kumar, V. Anomaly detection: A survey. *ACM computing surveys (CSUR)*, 41(3), 1-58, (2009).
10. Le, T. T. H., Kim, H., Kang, H., & Kim, H. (2022). Classification and explanation for intrusion detection system based on ensemble trees and SHAP method. *Sensors*, 22(3), 1154.
11. Friedman, J. H.: Greedy function approximation: a gradient boosting machine. *Annals of statistics*, 1189-1232 (2001).
12. Dietterich, T. G.: Ensemble methods in machine learning. In *International workshop on multiple classifier systems* (pp. 1-15). Berlin, Heidelberg: Springer Berlin Heidelberg, (2000, June).
13. Breiman, L.: Random forests. *Machine learning*, 45, 5-32 (2001).
14. Le, T. T. H., Oktian, Y. E., & Kim, H. (2022). XGBoost for imbalanced multiclass classification-based industrial internet of things intrusion detection systems. *Sustainability*, 14(14), 8707.
15. Le, T.T.H, Wardhani, R. W., Putranto, D. S. C., Jo, U., & Kim, H. (2023). Toward Enhanced Attack Detection and Explanation in Intrusion Detection System-Based IoT Environment Data. *IEEE Access*, 11, 131661-131676.
16. Recio-García, J. A., Orozco-del-Castillo, M. G., & Soladrero, J. A. (2023, July). Case-Based Explanation of Classification Models for the Detection of SQL Injection Attacks. In *ICCBR Workshops* (pp. 200-215).
17. Cumi-Guzman, B. A., Espinosa-Chim, A. D., Orozco-del-Castillo, M. G., & Recio-García, J. A. (2024). Counterfactual Explanation of a Classification Model for Detecting SQL Injection Attacks.
18. Le, T. T. H., Prihatno, A. T., Oktian, Y. E., Kang, H., & Kim, H. (2023). Exploring local explanation of practical industrial AI applications: A systematic literature review. *Applied Sciences*, 13(9), 5809.
19. Lundberg, S. (2017). A unified approach to interpreting model predictions. *arXiv preprint arXiv:1705.07874*.
20. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016, August). "Why should i trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 1135-1144).
21. Valeur, F., Mutz, D., & Vigna, G.: A Learning-Based Approach to the Detection of SQL Attacks. *Proceedings of the Conference on Detection of Intrusions and Malware & Vulnerability Assessment* (2005).
22. Gao, Q., Li, H., & Huang, T.: Detecting SQL Injection Attack Using an Artificial Neural Network. *Proceedings of the International Conference on Computational Intelligence and Security* (2008).
23. Xu, D., Huang, H., & Yin, H.: An Improved SQL Injection Detection Model Based on Machine Learning. *Journal of Software Engineering and Applications*, 3(12), 1131-1135. (2010).
24. Pan, X., Liu, L., & Yan, H.: SQL Injection Detection Based on AdaBoost Algorithm. *Proceedings of the IEEE International Conference on Computer and Information Technology* (2016).
25. Le, T. Q., Tran, D. H., & Nguyen, H. T.: SQL Injection Detection Using Gradient Boosting Decision Trees. *Proceedings of the International Conference on Information and Communication Technology* (2018).

26. Chen, T., & Guestrin, C.: XGBoost: A Scalable Tree Boosting System. Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (2016).
27. Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., & Liu, T. Y.: LightGBM: A Highly Efficient Gradient Boosting Decision Tree. Proceedings of the Advances in Neural Information Processing Systems (2017).
28. Nguyen, D. M., Tran, M. T., & Pham, B. T.: A Comparative Study of Machine Learning Algorithms for SQL Injection Detection. *Journal of Information Security and Applications*, 44, 144-154 (2019).
29. Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P.: SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, 16, 321-357 (2002).
30. Deriba, F. G., Salau, A. O., Mohammed, S. H., Kassa, T. M., & Demilie, W. B. (2022). Development of a compressive framework using machine learning approaches for SQL injection attacks. *Przeglad Elektrotechniczny*, 98(7), 181-187.
31. SQL Injection dataset, <https://www.kaggle.com/datasets/sajid576/sql-injection-dataset>, last accessed 2024/05/25.
32. Hosam, E., Hosny, H., Ashraf, W., & Kaseb, A. S. (2021, November). Sql injection detection using machine learning techniques. In 2021 8th International Conference on Soft Computing & Machine Intelligence (ISCMCI) (pp. 15-20). IEEE.
33. Gowtham, M., & Pramod, H. B. (2021). Semantic query-featured ensemble learning model for SQL-injection attack detection in IoT-ecosystems. *IEEE Transactions on Reliability*, 71(2), 1057-1074.
34. Abdulhamza, F. R., & Al-Janabi, R. J. S. (2022, November). SQL injection detection using 2D-convolutional neural networks (2D-CNN). In 2022 International Conference on Data Science and Intelligent Computing (ICDSIC) (pp. 212-217). IEEE.
35. Roy, P., Kumar, R., & Rani, P. (2022, May). SQL injection attack detection by machine learning classifier. In 2022 International Conference on Applied Artificial Intelligence and Computing (ICAIC) (pp. 394-400). IEEE.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.