

An Explainable-AI approach for Diagnosis of COVID-19 using MALDI-ToF Mass Spectrometry

Venkata Devesh Reddy Seethi, Zane LaCasse, Prajkta Chivte, Elizabeth R. Gaillard, Pratoool Bharti

Abstract—The novel severe acute respiratory syndrome coronavirus type-2 (SARS-CoV-2) caused a global pandemic that has taken more than 4.5 million lives and severely affected the global economy. To curb the spread of the virus, an accurate, cost-effective, and quick testing for large populations is exceedingly important in order to identify, isolate, and treat infected people. Current testing methods commonly use PCR (Polymerase Chain Reaction) based equipment that have limitations on throughput, cost-effectiveness, and simplicity of procedure which creates a compelling need for developing additional coronavirus disease-2019 (COVID-19) testing mechanisms, that are highly sensitive, rapid, trustworthy, and convenient to use by the public. We propose a COVID-19 testing method using artificial intelligence (AI) techniques on MALDI-ToF (matrix-assisted laser desorption/ionization time-of-flight) data extracted from 152 human gargle samples (60 COVID-19 positive tests and 92 COVID-19 negative tests). Our AI-based approach leverages explainable-AI (X-AI) methods to explain the decision rules behind the predictive algorithm both on a local (per-sample) and global (all-samples) basis to make the AI model more trustworthy. Finally, we evaluated our proposed method using a 70%-30% train-test-split strategy and achieved a training accuracy of 86.79% and a testing accuracy of 91.30%.

Index Terms—COVID-19 testing, Explainable-AI, Machine learning, RT-PCR test, Mass spectrometry, MALDI-ToF

I. INTRODUCTION

The virulent, fast spreading, and mutating nature of SARS-CoV-2 caused a worldwide pandemic with more than 229 million people infected and 4.7 million lives taken in the span of 1.5 years since emerging [1]. This pandemic has challenged the ability of healthcare systems to triage patients and overwhelmed capacities. Beyond healthcare, societal impacts have been apparent and major sectors of the global economy have been affected [2]. With this and the devastating nature of the pandemic, curbing the infection caused by the virus has become a major concern for all nations. Several protocols have been implemented such as mask mandates and lockdowns to reduce the spread of the virus. Frequent and rapid surveillance testing, however remains the most contributing factor in reducing the transmission of the virus [3]. The current method-of-choice for COVID-19 testing is based on reverse transcriptase-PCR (RT-PCR) platforms. These RT-PCR-based test procedures can be expensive, use complex reagents, and require a relatively high level of experience [4].

Alternatively, MALDI-ToF mass spectrometry (MS) methods provide advantages over RT-PCR-based methods such as the ability to rapidly process a large number of samples and yield highly-accurate results in COVID-19 diagnostic tests [5], [6]. The adoption of MALDI-ToF MS approach also sees a reduced demand for complex reagents like those required in molecular assay techniques such as RT-PCR. Apart from the equipment used, the specimen used in the analysis also plays a vital role in the applicability of a diagnostic

test. Current procedures use a nasopharyngeal (NP) swab or a saliva sample. Among these two specimens, saliva has been reported to contain higher viral loads for longer periods of time [7]. Moreover, saliva collection is less invasive and does not incite a cough, sneeze or a gag reflex, which could produce infectious aerosols [8].

For these reasons, in order to add to the number of diagnostic tests for COVID-19, we adopt the collection of saliva through gargle samples and leverage a MALDI-ToF MS-based approach for the analysis. This approach has already been established as an effective workflow, reported to agree with results from NP swab RT-PCR analysis for the diagnosis of COVID-19 [9], [10]. The output of such MALDI-ToF MS-based approach is a mass spectrum, which plots the mass-to-charge ratio (m/z) versus ion intensity as shown in Figure 1. This reflects the masses of biomolecules found within the saliva specimens and their relative abundance. These spectra can be complex, containing hundreds of signals that differ from specimen to specimen, and could include critical information regarding the presence of virus or the immune response of the patient.

Therefore, we adopt an AI-based approach to learn the relationships between the intensity of m/z peaks and the outcome of a RT-PCR test. Often, when an AI system learns a complex pattern, their decision making process is unclear and termed a black-box system. These black-box systems are uninterpretable and pose an impossible task for domain experts and health practitioners to understand the contributing factors to the outcome of a test. We tackle this problem using X-AI to present the rationale behind decision making of our AI-based approach. The X-AI approach not only provides valuable diagnostic information but also shows the ability of reproducing our results in clinical settings. The ability to analyze the response from the model instills trust in our AI system.

II. LITERATURE REVIEW

Following the recent success of AI solutions in the healthcare domain, several researchers have leveraged AI techniques to design efficient COVID-19 testing methods using mass spectrometry datasets [5], [6], [11]. These three studies were conducted on different datasets but the mode of sample collection, in this case NP swab, remains consistent. These studies differ, however, in sample processing and choice of algorithm. Nachtigall et al. [6] have used Cary-Blair transport medium to collect 362 samples, extracted and utilized peaks from the mass spectra as features. Further, they applied correlation-based [12] and information-gain-based [13] feature selection methods to identify 88 important peaks in the spectra. Although, they achieved 93.9% accuracy with the support vector machine (SVM) [14] algorithm, their results were not consistent across the population. In other work, Tran et al. [5] used a saline solution to collect 199 samples. While optimizing their AI solution, they used Machine Intelligence Learning Optimizer (MILO) [15] for hyperparameter searching. As a result, their extreme gradient boosting (XgBoost) [16] and deep neural network (DNN) [17] algorithms achieved 98.3% accuracy. In a more recent study, Deulofeu et al. [11] used nasal swabs and collected 148 samples using DeltaSwab-ViCUM viral transport medium (VTM-1) and 66 samples using DeltaSwab-Virus viral transport medium (VTM-2). Deulofeu et al. [11] pre-processed raw data

Venkata Devesh Reddy Seethi and Pratoool Bharti are with the Department of Computer Science at Northern Illinois University, Dekalb, Illinois, 60115 USA (email: devesh@niu.edu, pbharti@niu.edu)

Zane LaCasse, Prajkta Chivte, and Elizabeth R. Gaillard are with the Department of Chemistry and Biochemistry at Northern Illinois University, Dekalb, Illinois, 60115 USA (email: zlacasse@niu.edu, pchivte@niu.edu, gaillard@niu.edu)

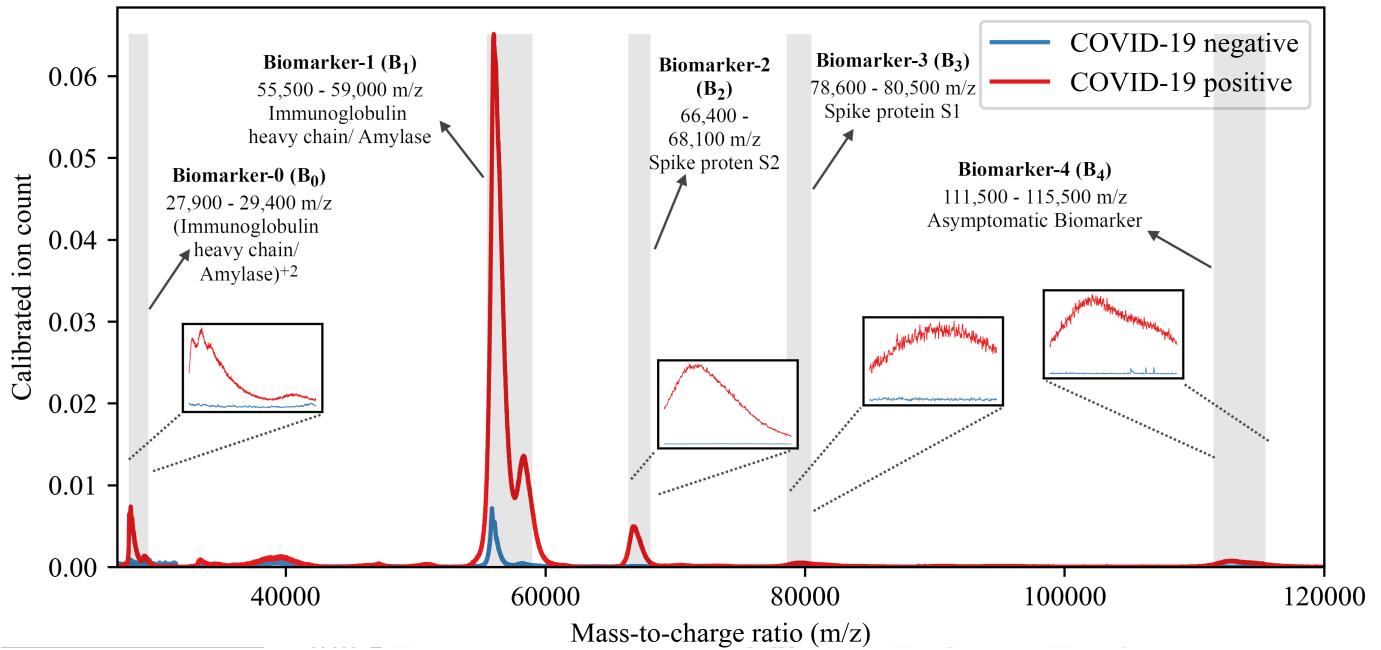


Fig. 1: Comparison of COVID-19 negative and a COVID-19 positive gargle MALDI-ToF spectra after preprocessing with baseline correction and normalization. Biomarkers ranges of interest in the spectra are highlighted with a gray hue. Zoom-in panels for each range show differences between the disease statuses.

using principal component analysis (PCA) [18] and evaluated several machine learning (ML) algorithms with a 10-fold cross-validation. The authors reported the accuracies for VTM-1 and VTM-2 datasets as 70.3% with XgBoost [16] and 88.4% with SVM [14] respectively.

The AI algorithms employed in these studies are complex and often require several pre- and post-processing steps. The absence of these details make it difficult to compare their results. Additionally, these solutions only provide a simple, binary decision: positive or negative. In order for medical practitioners and domain experts to further analyze their decisions and make more informed conclusions, it is crucial for testing methods to provide rationale on its decision making process. This ultimately makes the solution and decision-making process more trustworthy [19] for its application in clinical settings and more acceptable to end users. X-AI is one of the more recent advancements in AI [20] and is capable of accurately approximating the prediction mechanisms of AI algorithms. Explainable algorithms such as Shapely Additive Explanations (SHAP) [21] have gained popularity in the AI interpretability domain and AI in healthcare [22] but have never been used in COVID-19 diagnostic applications. In this paper, our goal is two fold: 1) to train an AI model to achieve high accuracy in COVID-19 diagnosis and 2) to leverage X-AI techniques to explain the reasoning behind the outcomes from AI models.

III. OUR METHODOLOGY AND CONTRIBUTIONS

The facts mentioned in section I emphasize the need for additional COVID-19 diagnostic methods that are relatively inexpensive, non-invasive, highly accurate, and are interpretable. Several studies since the onset of the pandemic have proposed more platforms and methods for COVID-19 testing, however, these have limitations based on the experiment procedures or the analysis. In this work, we address some of these limitations by interpreting our AI algorithm to provide more informed explanations of diagnosis, leveraging a cost-effective and easy-to-use MALDI-ToF MS approach, and using a non-invasive specimen type (saliva). After assessing multiple AI approaches, we

built our AI solution by training a tree-based ensemble algorithm (random forest [23]) on MALDI-ToF MS spectra data. Additionally, we interpret the predictions from trained random forest model using X-AI techniques such as SHAP, impurity-based feature importance, and permutation feature importance.

IV. DATASET PREPARATION

A. Data Collection

We collected gargle samples from students at Northern Illinois University Athletics (Illinois, USA) in the months of August and September of 2020. Students who consented to participate in our study were asked to provide a water gargle sample and a simultaneous NP swab sample. The gargle sample was processed as described in [10] and the NP swabs were confirmed COVID-19 positive/negative by RT-PCR conducted by the Dekalb County Health Department (Illinois, USA). For each gargle sample, the outcome of the RT-PCR test on NP swab was used as the ground truth in our analyses. In total, 152 gargle samples were analyzed, including 60 COVID-19 positive and 92 COVID-19 negative.

In brief, sample handling included participants gargling 10 mL of bottled spring water for 30 seconds which was then deposited in a 50 mL conical centrifuge tube. The collected specimen was filtered through a $\leq 0.45\mu\text{m}$ syringe filter, acetone precipitated to concentrate proteins, centrifuged, and resuspended in a reconstitution buffer. After incubation, we spotted each sample on a target plate, air-dried, and transferred the plate to a Shimadzu AXIMA Performance [24] mass spectrometer, which was calibrated with a protein standard daily. Ions were detected in positive-ion, linear detection mode over a range of 2,000 – 200,000 m/z. Finally, we corrected the baseline of all the generated spectra using the Shimadzu Launchpad software obtained from the manufacturer [24] and normalized the ion count intensities with the parent peak of the protein standard calibrant.

B. Potential Biomarker Selection

The presence or absence of peaks and their relative intensities in the collected mass spectra can be correlated to the presence or abundance of proteins within each specimen. The spectra we collected represent a multitude of proteins (after acetone precipitation) that are in the range of 2,000 – 200,000 m/z. Since our goal was to differentiate between a positive and negative sample, we carefully selected only those m/z ranges that were observed to change given the COVID-19 status by RT-PCR. An overlay of mass spectra highlighting these differences for a COVID-19 positive and negative sample is shown in Figure 1. Note that Figure 1 only shows a segment of the spectra where potential COVID-19 biomarkers are present. An attempt to identify the potential biomarkers was made in our previous work [10]. These biomarker peaks were selected to potentially include both human proteins (immunoglobulin heavy chain, amylase, asymptomatic biomarker) and SARS-CoV-2 viral proteins (Spike protein subunits S2 and S1), summarized in Table I. For the sake of simplicity, here, these ranges were named Human Protein-1, Human Protein-2, Human Protein-3, and Viral Protein-1, Viral Protein-2, respectively.

Biomarker	Mass Range	Potential Identity
B_0	27,900 – 29,400	Human Protein-1
B_1	55,500 – 59,000	Human Protein-2
B_2	66,400 – 68,100	Viral Protein-1
B_3	78,600 – 80,500	Viral Protein-2
B_4	111,500 – 115,500	Human Protein-3

TABLE I: Potential biomarker ranges and their descriptions.

C. Feature Engineering

The performance of supervised ML algorithms depends on the quality of the features used in training. Typically, the prediction power of any ML classification algorithm increases with an increase in the number of features and languishes after reaching a certain number [25]. If we had directly used the raw readings of all five chosen m/z ranges within the spectra, each sample would have 12,600 features (the combined bandwidth of all 5 ranges); this could potentially overfit the ML models for our limited dataset of 152 samples. We thus extracted features that balanced being easy to compute but also have enough information to discriminate between the positive and negative samples. After iterating through many statistical and intuitive features, we decided to compare the area under the curves (AUC) [26] for the five biomarkers ($B_0 - B_4$). Although the AUC captures the intensity of a biomarker, a wide variance from one sample to another was observed. This happens mainly due to two reasons. First, each sample has potentially different amounts of human/viral protein present and second, information regarding exposure to the virus or time course of infection was not collected. Such differences in immune response and viral load would be reflected in the biomarker intensities. Regardless, we observed that the ratio of AUC between pairs of biomarkers showed a consistent pattern for positive and negative samples. Therefore, we computed 10 features for each sample as the ratio of each possible combination of different biomarker pairs (shown in Table II). These features were computed for the MALDI-ToF dataset of size 152 with 60 COVID-19 positive and 92 COVID-19 negative samples.

D. Data Preparation

We partitioned 152 samples from our MALDI-ToF dataset using a train-test-split strategy with 70 : 30 train-test-split ratio to form two disjoint and independent splits. This prevents data leak from classifier training into performance assessment. Within both train- and

test-splits, the ratio of the count of positive to negative samples was kept constant to maintain class balance. Our classifier uses the train-split dataset (41 COVID-19 positive and 65 COVID-19 negative) and learns the COVID-19 prediction task. We then assess the performance of the trained classifier on the unseen test-split dataset (19 COVID-19 positive and 27 COVID-19 negative samples).

V. METHODS

After evaluating several tree-based algorithms, we selected random forest (RF) [23] classifier due to its superior performance among tree-based algorithms. The goal of RF algorithm is to train a classifier function (f) that learns the relationships between the ratio features generated from training samples, (x_{train}) and the outcome of the associated COVID-19 RT-PCR test (y_{train}). After completion of training, the algorithm was tested on new/unseen saliva samples (x_{test}) to evaluate the precision and recall metrics of the trained classifier ($f(x_{test}) = y_{test}$) as compared to the corresponding RT-PCR result.

Further, we present the complex relationships learned by the trained RF algorithm between the input ratio features and the outcome of a COVID-19 RT-PCR test in a human interpretable manner with the aid of X-AI algorithms. By leveraging the X-AI algorithms, explanations for the outcomes of RF algorithm are possible on both a global (all-samples) and local (per-sample) basis.

Feature	Description
R_0	AUC_{B_0} / AUC_{B_1}
R_1	AUC_{B_0} / AUC_{B_2}
R_2	AUC_{B_0} / AUC_{B_3}
R_3	AUC_{B_0} / AUC_{B_4}
R_4	AUC_{B_1} / AUC_{B_2}
R_5	AUC_{B_1} / AUC_{B_3}
R_6	AUC_{B_1} / AUC_{B_4}
R_7	AUC_{B_2} / AUC_{B_3}
R_8	AUC_{B_2} / AUC_{B_4}
R_9	AUC_{B_3} / AUC_{B_4}

TABLE II: Ratio features generated from AUC values of five potential biomarkers.

A. Random Forest (RF) Classifier

RF is an ensemble classification algorithm that combines several tree predictors (Decision Trees [27]) to build a forest classifier. Each tree in the forest is trained from a randomly sampled vector from the input features. These random vectors are sampled independently from the same distribution in the predictor space [23]. The outcome of a trained forest is the majority vote taken in the form of mean or median of the predictions from all the independently trained trees. Additionally, in RF algorithm, we can fine-tune the parameters such as the number of trees (also known as estimators) in the forest, methods of bootstrapping, and the properties of trees. If bootstrapping is enabled, random samples are taken from the training set with replacement to train decision trees and take the majority vote. We can tune the core properties of each decision tree such as maximum depth, maximum number of leaves, splitting criterion, minimum samples per split, and minimum samples per leaf to optimize the performance.

The first step in training the RF classifier is to find optimized hyperparameters by leveraging a grid search technique that uses 5-fold cross-validation on the training dataset. The grid search algorithm conducts exhaustive try-outs on all possible sets of hyperparameter configurations wherein each set of configuration is catalogued with its cross-validation accuracy score. The hyperparameter configuration with the best score (shown in Table III) was selected to train the

RF algorithm on full training data which is further evaluated and interpreted with the test-split data.

Hyperparameter	Value
Number of estimators	100
Bootstrap	True
Splitting criterion	Gini impurity
Maximum depth of tree	3
Minimum samples per split	2
Minimum samples per leaf	1

TABLE III: Hyperparameters of our RF classifier.

B. Random Forest Interpretation Using Explainable-AI

To instill trust in the prediction of the model, we interpret the RF model both globally and locally using X-AI techniques. While global interpretation explains the importance of each feature in overall model prediction, local interpretation provides the reasoning on a single instance level. For global interpretation, analyzing the importance of features determined by a trained model gives information about its prediction process. In this section, we leverage multiple techniques to investigate feature importance of our trained RF model.

1) *Impurity-based feature importance (IFI)*: First, we interpreted the trained RF classifier with a model-specific, IFI technique. This technique computes the feature importance as the average decrease in Gini impurity [28] from a feature node to its child nodes. Generally, except for the leaf nodes that occur at the bottom of the tree, all other nodes have a threshold condition placed on a feature whose decision splits the node into two child nodes. The drop in impurity in this split is then calculated as the difference of Gini impurity between the node and the sum of its child nodes, weighted by the number of samples in each node. The final feature importance is the mean of drop in impurity at every prevailing feature node across the forest. Following these procedures, we calculated the IFI for all 10 features as shown in panel-A of Figure 2 with the drop in Gini impurity on the y-axis and features on the x-axis. Although this method gives a general understanding of feature rankings, one disadvantage is inflation of importance for features with high cardinality in the forest. Therefore, we investigated other feature importance algorithms to validate the results.

2) *Permutation Feature Importance (PFI)*: The PFI technique [23] calculates the drop in performance when the trained RF algorithm is validated with permuted feature values which essentially makes the feature a noise. This method analyzes one feature at a time by permuting their values to break the association between the feature and the outcome. Clearly, permutation of important features results in a larger performance drop. Therefore, the feature importance is gauged as the drop in permuted feature accuracy from the base accuracy (where every feature was used in its original form). Bias from random value assignment to feature values was prevented by permuting values and averaging the drop in accuracy across 100 iterations. However, this algorithm falls short in capturing feature interaction effects as it permutes only one feature at a time. Panel-B in Figure 2 shows the permutation importance scores for RF classifier with overall feature importance printed on the y-axis and the features on the x-axis.

3) *Shapely Additive Explanations (SHAP)*: SHAP algorithm, introduced by Lundberg and Su-In [21] can be formalized as an additive feature attribution technique that unifies concepts from local model-agnostic explanations (LIME) [19] and Shapely values [29]. In general, linear models are easy to interpret as the feature importance of the model are reflected by the coefficients of the model itself. This is not the case for the RF algorithm as a simple linear model

does not have the capability to approximate a more complex RF algorithm. However, output of a single sample could be explained by generating some neighbourhood samples followed by fitting a linear model locally on the newly generated neighborhood samples. The LIME algorithm follows this approach: it generates a local point of view per sample by applying perturbations to its feature values and generating some artificial data in the neighbourhood of a sample. LIME then fits a surrogate model on the artificial data and presents the coefficients of the surrogate model as the model interpretation. The downside of LIME lies in its generation of artificial samples using a heuristic kernel function. Artificially replicating clinical data is a difficult task and depending on kernel function heuristics remains unreliable. SHAP adopts LIME procedures and bridges its limitations with the Shapely values which borrow principles from cooperative game theory. Shapely values rationalize the distribution of credit or payoff among the players in a cooperative game by inspecting the contribution of each player individually as well as in a collaborative team effort [29].

In theory, the feature attributions assigned by our trained RF classifier (f) to the features R_i (where $i \in (0,9)$ for 10 ratio features) in a given datasample (x) is computed as $(\phi_{R_i}(f, x))$ using Equation (3). This equation facilitates SHAP to compute a unique optimal solution for feature attribution by combining individual as well as team efforts. This can be understood in light of PFI where we extracted feature attributions by masking each feature at a time and computing the drop in performance of the classifier. Clearly, PFI captures individual contributions but fails to grasp collaborative contributions. The more intuitive approach of SHAP allows for collaborative contributions in conjunction to individual contributions to be considered by iterating over all possible subsets S in our feature space having F number of features excluding feature R_i (where $F = 9$) as $S \subseteq F/\{R_i\}$. In each subset, the marginal contribution (C) of feature R_i is extracted as the difference in outcome of function f with and without R_i in the subset as shown in Equation (1) where $f_{S \cup \{R_i\}}$ and f_S are two retrained functions with and without R_i respectively that are marginalized over features absent in the subset. The resultant sum of C is weighted with W (refer to Equation (2)) whose numerator is the product of different number of ways a subset can be formed ($S!$) and the number of ways to choose features excluding the features in a subset $((|F| - |S| - 1)!)!$. The numerator is normalized by $F!$, the total possible ways of choosing the features. Finally, for each feature in x , SHAP values are computed using Equation 3 that iterates over all possible subsets S and calculates SHAP values. The final SHAP values have a magnitude that shows the importance of the feature and a direction where positive sign or negative sign corresponds to their contribution towards the COVID-19 positive or negative decision.

Since the idea behind Shapely values is that there exists only one solution to allocate the attributions, the results are more reliable than previous methods such as LIME [19]. In addition to this, SHAP has other desirable axiomatic properties. The local accuracy property entails that the sum of all feature attributions of a model should approximate the original model meaning $f(x) = \phi_{R_n} + \sum_{i=0}^9 \phi_{R_i}$, where R_n is the feature attribution when all features are toggled off. The missingness property ensures that the features that are toggled off have no impact in model's output. Lastly, the consistency property ensures for any feature R_i , ϕ_{R_i} remains consistent with respect to the impact of feature values of on any retrained model. For example, in two functions trained on different subsets including the feature R_i , if the feature values of R_i have larger impact on $f1$ than $f2$, then $\phi_{R_i}(f1, x) \geq \phi_{R_i}(f2, x)$.

$$C = f_{S \cup \{R_i\}}(x_{S \cup \{R_i\}}) - f_S(x_S) \quad (1)$$

$$W = \frac{|S|!(|F| - |S| - 1)!}{|F|!} \quad (2)$$

$$\phi_{R_i}(f, x) = \sum_{S \subseteq F / \{R_i\}} W \times C \quad (3)$$

One can assume that classical Shapely value approach in Equation (3) is a promising way to obtain feature attributions, nevertheless, the computational complexity reaches non-polynomial time which makes it difficult to apply to our dataset. In 2020, Lundberg [22] introduced tree explainer, a model that leverages the underlying structure of tree-based models and optimizes the computational complexity to polynomial time. Tree explainer approximates the model while adhering to the natural properties of Shapely values [29] mentioned previously. Therefore, we used tree explainer [22] to generate a SHAP value matrix of dimensions [46, 10], for 46 testing samples having 10 ratio features. Using the SHAP value matrix, we then interpret our trained RF model both globally and locally. For global model interpretation, we extract the magnitude of feature importances in two steps: 1) take the absolute SHAP values which nullifies the effect of direction and 2) compute the mean of absolute SHAP values across all samples. This results in condensing the [46, 10] matrix to a single array of 10 feature importances which are shown in panel-C of Figure 2. By comparing all 3 feature importance algorithms in Table IV, it is conclusive that the two most important features are R_9 and R_5 and the least important are R_2 and R_4 . In addition to global (all-sample) explanations, we present local (per-sample) explanations using SHAP value matrix as shown in Figure 6. In these plots, the x-axis represents the SHAP value per-sample and features are numbered on the y-axis. Additionally, we embedded feature values in color ranging from blue (low value) to red (high value) and assigned each feature to a quartile based on its value as shown at origin of each bar. This enables us to interpret the SHAP values and observe their effects on the feature values. Details on global and local explanation are provided in sections VI-B and VI-C, respectively.

VI. RESULTS

A. Results of Random Forest Classifier

We trained the RF algorithm with the hyperparameters (shown in Table III) extracted through an exhaustive grid search. A total of 106 samples were employed in training and 46 in testing the model. Although data augmentation techniques could enrich the dataset to create more training samples, we avoided doing this because it could have made the dataset noisy by adding false patterns. The trained RF algorithm achieved 91.30% accuracy on the test-split, whereas 86.79% accuracy was achieved on the train-split. Similar high accuracies on both training and test sets assure that the model is not overfitted. Confusion matrix for both splits are shown in Figure 3. Overall, only 4 and 14 samples are classified falsely in the test and train dataset, respectively. As mentioned earlier, our dataset had more negative samples than positive ones, making it crucial to examine the accuracy of the model for both positive and negative categories. Therefore, we evaluated precision, recall and F1 metrics on both the test and train dataset for positive and negative samples independently. As shown in the evaluation report (in Figure 4), the train-split had a high precision, recall and F1-score of 0.86, 0.94, 0.90 for the COVID-19 negative samples and 0.89, 0.76, 0.82 for positive samples, respectively. On the other hand, precision, recall and F1-scores for COVID-19 negative and positive for the test-split were 0.87, 1.00, 0.93 and 1.00, 0.79, 0.88, respectively. In case of test-split, there were 0 false positives, 4 false negatives and other 42 samples had true predictions. In train-split, the false predictions were

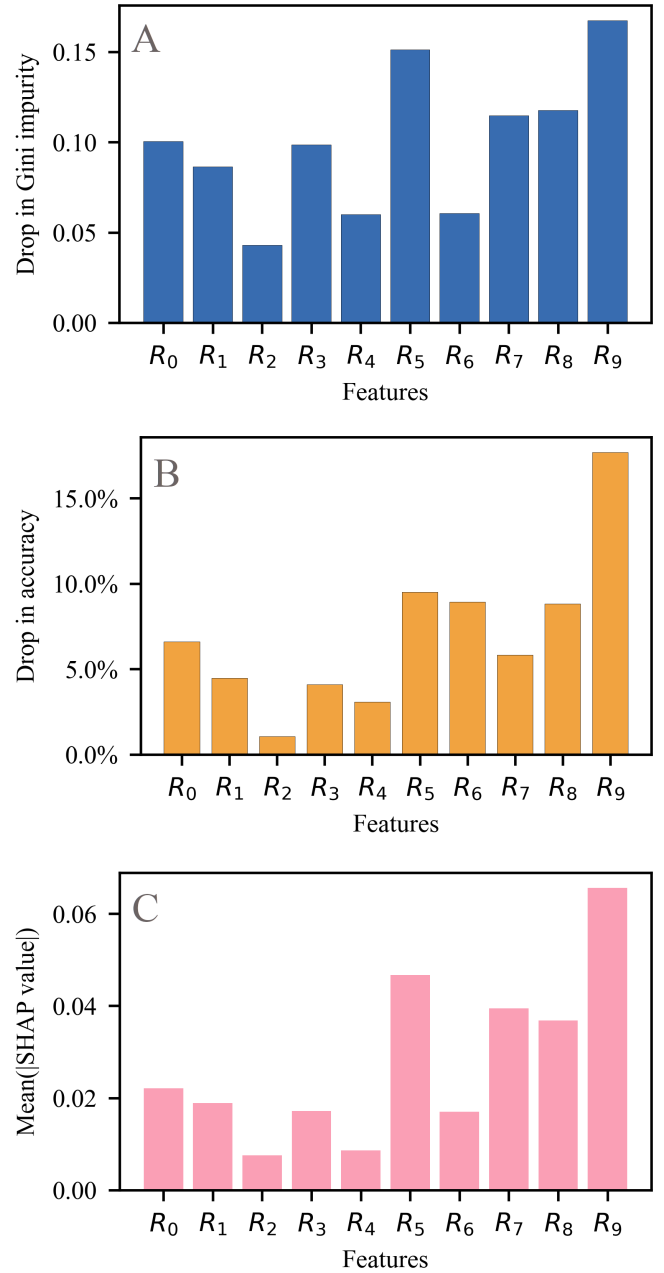


Fig. 2: Feature importance for RF classifier using three feature interpretation techniques: Impurity-based feature importance (panel-A, top), Permutation Feature Importance (panel-B, middle), and Shapely Additive Explanations (panel-C, bottom). The y-axis shows the feature importance for each feature shown on the x-axis.

slightly higher with 4 false positives and 10 false negatives and a total of 92 correct predictions. Although these metrics portray high efficiency of the trained model in COVID-19 testing, it still does not provide transparency on its prediction mechanisms. This motivated us to further dissect the trained RF model and understand its prediction mechanisms using X-AI techniques.

B. Global Interpretation of Random Forest Classifier

The idea of global interpretation of a model goes beyond the analysis of feature importance. In fact, it can elaborate on how the distribution of each feature is affecting the model prediction. As shown earlier, each of the three feature importance algorithms concluded

Method	Order of feature importances
Impurity-based feature importance	$R_9 > R_5 > R_8 > R_7 > R_0 > R_3 > R_1 > R_6 > R_4 > R_2$
Permutation feature importance	$R_9 > R_5 > R_6 > R_8 > R_0 > R_7 > R_1 > R_3 > R_4 > R_2$
Shapely additive explanations	$R_9 > R_5 > R_8 > R_7 > R_0 > R_6 > R_3 > R_1 > R_4 > R_2$

TABLE IV: Global feature importances shown in descending order for different feature interpretation techniques.

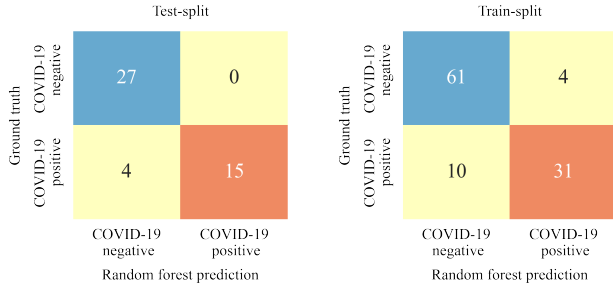


Fig. 3: Confusion matrices for the RF classifier on the train-split and test-split.

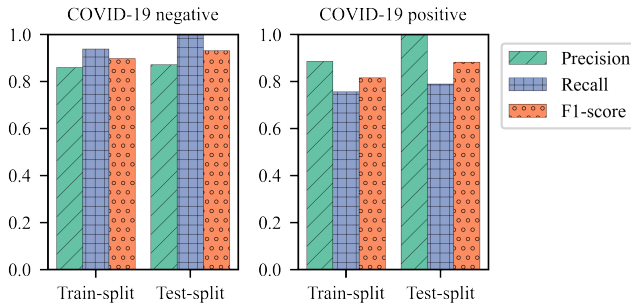


Fig. 4: Precision, recall and F1-score for test-split dataset for COVID-19 negative (on left) and COVID-19 positive (on right).

similar feature ranking: the top two features were consistently R_9 and R_5 while the bottom two were R_2 and R_4 (Table IV). On the other hand, R_7 , R_8 , R_1 , and R_3 had slightly varying importances among the three algorithms. Only one feature, R_6 , was seen to be varying by two positions across them. It is fairly common to observe slight differences in the feature importances across the three algorithms due to underlying procedural differences. However, on the whole, these rankings reveal that nine out of ten feature importances assigned by the three algorithms were fairly consistent. The consistent nature of feature attributions show that our algorithm can be understood when viewed from the different perspectives of the X-AI algorithms. Although these are important findings, it doesn't explain how the model prediction is correlated with the feature values. While IFI and PFI algorithms are good at generating feature rankings, they do not provide enough information to explore the relationship between model prediction and the distribution of individual feature. Therefore, we extended the SHAP algorithm that is capable to analyze the feature distribution and the model dependency on it.

The SHAP algorithm generated a violin plot, which is a fine-grained view of SHAP (feature importance) values and feature distribution as shown in Figure 5. This plot illustrates the distribution of SHAP values across all test-split samples, allowing us to observe trends of SHAP values in conjunction with feature values. To interpret these plots, the higher the SHAP value, the larger is the importance; the positive SHAP values support positive COVID-19 prediction,

negative values support negative output; and the feature value is color coded where blue represents the values in lower ranges, red represents the values in higher ranges. For example, the feature R_9 has SHAP values extending up to $+0.165$ towards a positive COVID-19 decision, where the feature value is low as it is colored in blue. On the negative side, it can be seen that mostly high values of R_9 contribute towards a negative COVID-19 prediction. On the contrary for feature R_5 , smaller and larger feature values contribute towards COVID-19 negative and positive tests, respectively. Features R_6 and R_7 follow similar patterns as R_5 while R_0 and R_1 have similar patterns of R_9 .

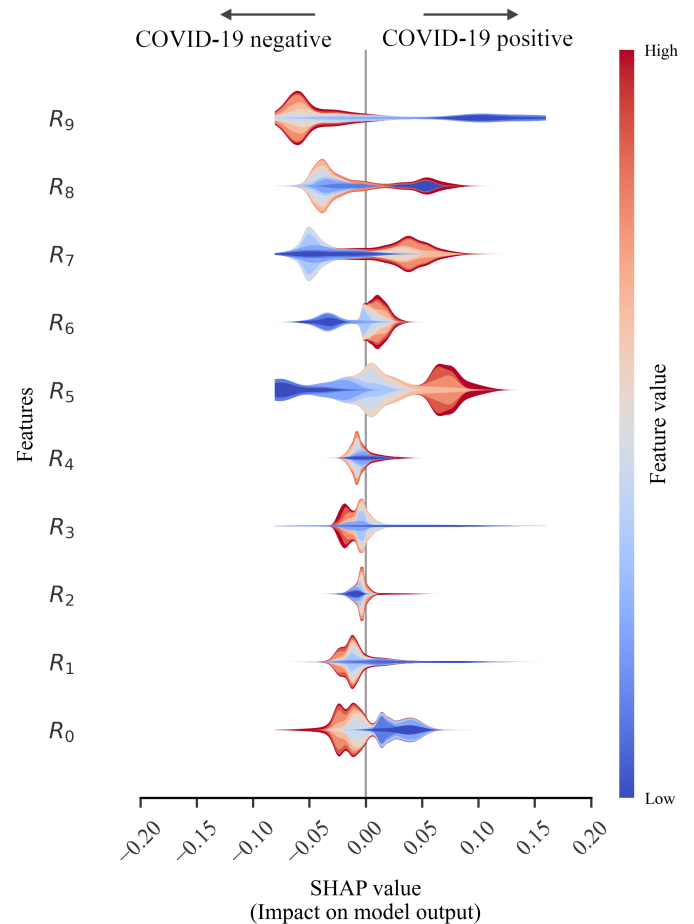


Fig. 5: Violin plot showing for global interpretation of RF classifier using SHAP values.

The violin plot provides an approximate idea about the model output for any sample based on the feature distribution. This can be immensely helpful to cross examine the model in the case that the model prediction is not following the trend of the violin plot. For example, if the model is predicting COVID-19 positive for a sample where the R_9 feature value is very high, we can easily identify that the model is not following the trend and hence, it needs further validation. This kind of information is critical to instill trust in users so that they can question outputs from the model.

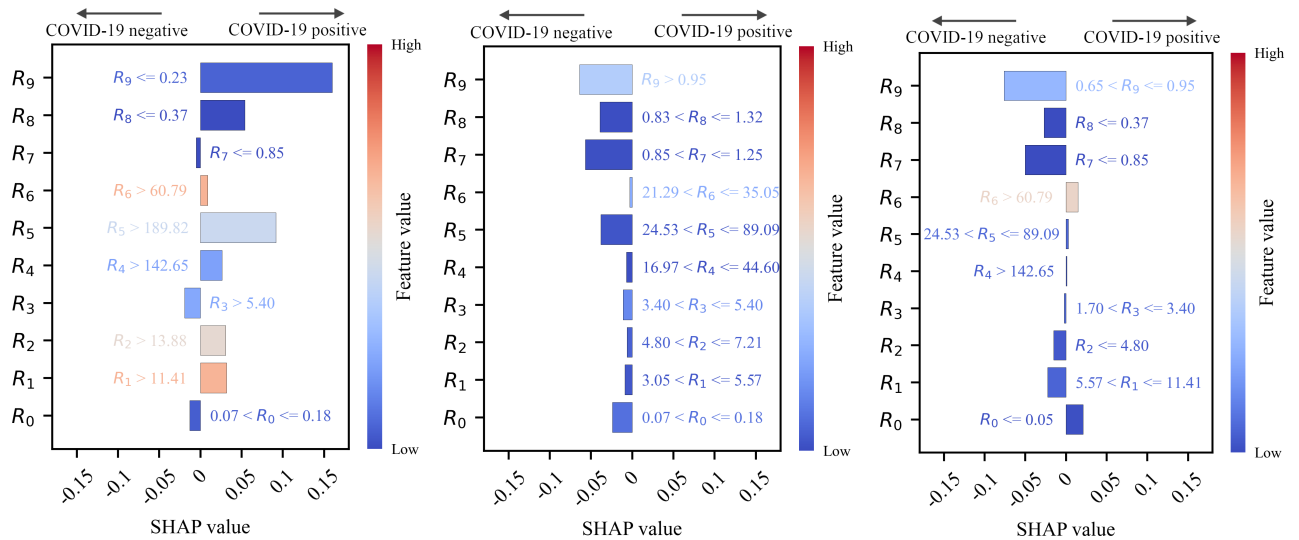


Fig. 6: Local explanations generated using SHAP algorithm for three records, true positive sample (on the left), true negative sample (in the middle) and false negative sample (on the right).

C. Local Interpretation of Random Forest Classifier

For domain experts and clinicians handling COVID-19 testing, local explanations of the model unveil rich diagnostic information as they provide rationale behind predictions on a per-sample basis. For this purpose, we have further leveraged the SHAP algorithm for interpreting one sample at a time. To do this, we take a holistic view on the classifier performance by viewing different cases of predictions such as a true positive (TP), a true negative (TN), and a false negative (FN) (shown in Figure 6). Since, there were no cases of false positives in the test dataset, we could not explore them in examples. In Figure 6, each bar is colored based on its feature value, ranging from a small value (blue) to large value (red). The labels at the origin of each bar indicate the condition on each feature value to support either category. The feature importance, similar to the global SHAP plots, is shown on the x-axis.

The local interpretation helps us to understand how, for any individual sample, different features interact with each other to sway the final decision of the model. In the case of a true positive sample (left plot in Figure 6), R_9 and R_5 are the most contributing features towards COVID-19 positive outcome. followed by R_8 , R_4 , R_2 and R_1 with lower importance. Although R_0 , R_3 and R_7 are directing towards a negative prediction, there is more support for the positive case based on the SHAP values. In the true negative record (middle plot in Figure 6), the R_9 feature value is mid-ranged followed by features R_7 , R_8 , R_5 , and R_0 . We also observe in the true negative sample that all features except R_9 have relatively smaller values compared to a true positive. Although not any single feature is directing very strongly towards negative test, all of them together support a negative prediction. We finally show explanation of a false negative case (right plot in Figure 6). In the false negative sample, the R_9 feature had a smaller value compared to the true negative sample. However, all other features except R_6 had similar values compared to the true negative sample. These high similarities of COVID-19 positive sample with a negative one potentially alluded the model to classify a COVID-19 positive as a negative test. Perhaps a medical practitioner in control of making decision based on the outcome of a COVID-19 test can benefit from a local analysis of the test and inspect the features that contributed most towards a positive or negative test. This in conjunction with consideration of a patient's symptoms can lead to more informed conclusions and diagnosis.

VII. DISCUSSION

While analyzing our results through various interpretation techniques, we observed quite remarkably that all three methods of global feature interpretation consistently resulted in the most important features as R_9 and R_5 (Table IV). The meaning of these features is of particular interest because of the biomarker ranges making up the features, B_3/B_4 for R_9 and B_1/B_3 for R_5 . First, the features R_9 and R_5 are host-to-virus and virus-to-host features, respectively, and represent the interaction between host immune response and viral antigen. This is as described in our previous work [10]. Second, the viral protein fragment that corresponds to B_3 is tentatively suggested to be the S1 subunit of the viral spike protein (see Table I). S1 is a key viral biomarker in the mass spectra because the observed m/z ratio does not overlap significantly with other proteins or fragments from the gargle sample.

The least important features, R_4 and R_2 , were decidedly the same by each method as well. R_4 contains B_2 , which is ambiguous in that the AUC value for this peak may consist of *both* human and viral proteins (a small amount of human serum albumin and spike protein S2 subunit, respectively). Therefore, the presence of B_2 could render R_4 less important, as B_2 would be present in both disease and non-disease states. Considering the nature of the biomarkers, it was odd that R_2 was not found to be as important as R_5 : both ratio features should technically represent the same proteins, as B_0 is the double charged ion of B_1 . To investigate this, the correlation plot between each ratio feature value was plotted and, not surprisingly, R_5 and R_2 showed a linear relationship (as shown in Figure 7). This is most likely the reason why R_2 was deemed less important than R_5 and the least important feature altogether: the algorithm avoids repetitive predictions by not placing importance on features equivalent in meaning. It should be pointed out that each biomarker has been assigned a putative identity based on preliminary studies only and have yet to be unambiguously characterized.

Turning to local feature interpretation, these could offer a powerful diagnostic view: instead of a binary decision by RT-PCR or algorithm with a percent confidence, the local explanations generated using SHAP values show the contribution of each feature in each sample. This could allow domain experts and health officials administering such a test to make more informed decisions based on the putative host-to-viral or viral-to-host ratios, for case-by-case bases as seen

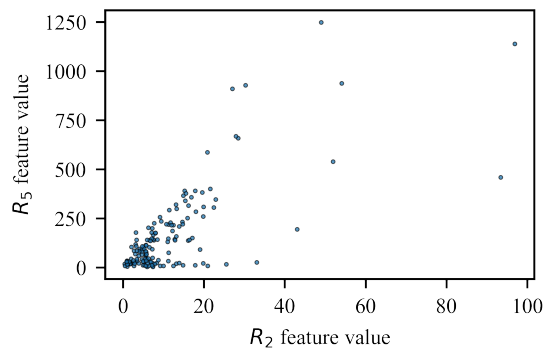


Fig. 7: Correlation plot showing positive correlation between feature values of R_5 and R_2 .

in Figure 6. If the algorithm decides positive or negative, the local feature explanation would provide additional information as to *why* the test output its decision. However, due to the small size of the dataset and the lack of information regarding time points of infection/exposure when the samples were collected, it is not possible to draw definitive conclusions from the local explanations presented here.

Nonetheless, not only in this diagnostic approach is the viral aspect of COVID-19 considered (as is only considered in PCR-based tests), but the human immune response is also given due consideration. ML algorithms coupled with the utility of X-AI interpretability can be implemented into healthcare to address the limitations of black box systems and improve the diagnostic process, from test output to informed care.

VIII. CONCLUSION

We have analyzed the application of machine learning models on the MALDI-ToF based data of gargle samples and inferred the performance of RF to be the best among the other algorithms. The RF classifier was evaluated on a 70%-30% train-test-split strategy where the accuracy on the test dataset was 91.30%. Our observations strongly support that high accuracy for COVID-19 diagnosis can also be achieved using saliva samples as compared to NP swabs. Given the ease-of-use and the advantages of saliva as a specimen over NP swabs, gargle samples can be employed in clinical settings. Further, we utilized the concepts of explainable-AI and interpreted the RF algorithm using SHAP and feature importance techniques such as permutation feature importance and impurity-based feature importance. By using these interpretation models, we showed both global and local explanations to the model's performance. This step is essential as it lets the domain experts and medical practitioners to understand the mechanisms of the black-box AI models when using in a real-time basis. Our approach will help developing faster, cheaper diagnostic methods for COVID-19 which are not only precise but also interpretable and trustworthy.

IX. ACKNOWLEDGEMENTS

- This research used computing resources of data, devices and interaction Laboratory (ddiLab) at Northern Illinois University.
- The authors gratefully acknowledge instrumentation support from Shimadzu Scientific Instruments, Inc.

REFERENCES

[1] "WHO COVID-19 Dashboard," <https://covid19.who.int/>.

- [2] C. Sohrabi, Z. Alsafi, N. O'Neill, M. Khan, A. Kerwan, A. Al-Jabir, C. Iosifidis, and R. Agha, "World health organization declares global emergency: A review of the 2019 novel coronavirus (COVID-19)," *International Journal of Surgery*, vol. 76, pp. 71–76, April 2020.
- [3] M. Nicola, C. Sohrabi, G. Mathew, A. Kerwan, A. Al-Jabir, M. Griffin, M. Agha, and R. Agha, "Health policy and leadership models during the COVID-19 pandemic-review article," *International Journal of Surgery*, vol. 81, p. 122–129, September 2020.
- [4] L. J. Carter, L. V. Garner, J. W. Smoot, Y. Li, Q. Zhou, C. J. Saveson, J. M. Sasso, A. C. Gregg, D. J. Soares, T. R. Beskid, S. R. Jervey, and C. Liu, "Assay techniques and test development for COVID-19 diagnosis," *ACS Central Science*, vol. 6, no. 5, pp. 591–605, April 2020.
- [5] N. K. Tran, T. Howard, R. Walsh, J. Pepper, J. Loegering, B. Phinney, M. R. Salemi, and H. H. Rashidi, "Novel application of automated machine learning with MALDI-ToF-MS for rapid high-throughput screening of COVID-19: a proof of concept," *Scientific Reports*, vol. 11, no. 1, pp. 1–10, April 2021.
- [6] F. M. Nachtigall, A. Pereira, O. S. Trofymchuk, and L. S. Santos, "Detection of SARS-CoV-2 in nasal swabs using MALDI-MS," *Nature Biotechnology*, vol. 38, no. 10, pp. 1168–1173, July 2020.
- [7] T. K. Burki, "Testing for COVID-19," *The Lancet Respiratory Medicine*, vol. 8, no. 7, pp. e63–e64, May 2020.
- [8] E. Pasomsab, S. P. Watcharananan, K. Boonyawat, P. Janchompoo, G. Wongtabtim, W. Suksuwan, S. Sungkanuparph, and A. Phuphuakrat, "Saliva sample as a non-invasive specimen for the diagnosis of coronavirus disease 2019: a cross-sectional study," *Clinical Microbiology and Infection*, vol. 27, no. 2, pp. 285.e1–285.e4, February 2021.
- [9] R. K. Iles, R. Zmuidinaite, J. K. Iles, G. Carnell, A. Sampson, and J. L. Heeney, "Development of a clinical MALDI-ToF mass spectrometry assay for SARS-CoV-2: rational design and multi-disciplinary team work," *Diagnostics*, vol. 10, no. 10, p. 746, August 2020.
- [10] P. Chivte, Z. LaCasse, V. D. R. Seethi, P. Bharti, J. Bland, S. S. Kadkol, and E. R. Gaillard, "MALDI-ToF protein profiling as a potential rapid diagnostic platform for COVID-19," *Journal of Mass Spectrometry and Advances in the Clinical Lab*, vol. 21, pp. 31–41, August 2021.
- [11] M. Deulofeu, E. García-Cuesta, E. M. Peña-Méndez, J. E. Conde, O. Jiménez-Romero, E. Verdú, M. T. Serrando, V. Salvadó, and P. Boadas-Vaello, "Detection of SARS-CoV-2 infection in human nasopharyngeal samples by combining MALDI-ToF-MS and artificial intelligence," *Frontiers in Medicine*, vol. 8, p. 398, April 2021.
- [12] M. A. Hall, "Correlation-based feature selection for discrete and numeric class machine learning," in *Proceedings of the Seventeenth International Conference on Machine Learning*, June 2000, p. 359–366.
- [13] S. Appavu, R. Rajaram, M. Nagammai, N. Priyanga, and S. Priyanka, "Bayes theorem and information gain based feature selection for maximizing the performance of classifiers," in *International Conference on Computer Science and Information Technology*, vol. 131, January 2011, pp. 501–511.
- [14] W. S. Noble, "What is a support vector machine?" *Nature Biotechnology*, vol. 24, no. 12, pp. 1565–1567, December 2006.
- [15] H. Rashidi, N. Tran, and S. Albahra, "Machine intelligence learning optimizer (MILO)," <https://milo-ml.com/>.
- [16] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '16. Association for Computing Machinery, August 2016, p. 785–794.
- [17] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, May 2015.
- [18] S. Wold, K. Esbensen, and P. Geladi, "Principal component analysis," *Chemometrics and Intelligent Laboratory Systems*, vol. 2, no. 1–3, pp. 37–52, August 1987.
- [19] M. T. Ribeiro, S. Singh, and C. Guestrin, "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*, August 2016, p. 1135–1144.
- [20] P. Linardatos, V. Papastefanopoulos, and S. Kotsiantis, "Explainable AI: A review of machine learning interpretability methods," *Entropy*, vol. 23, no. 1, p. 18, January 2021.
- [21] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, December 2017, p. 4768–4777.
- [22] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, and S.-I. Lee, "From local explanations to global understanding with explainable AI for trees," *Nature Machine Intelligence*, vol. 2, no. 1, pp. 56–67, January 2020.
- [23] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, October 2001.

- [24] Shimadzu, "Axima performance - a highly flexible research-grade mass spectrometer," <https://www.shimadzu.com/an/products/maldi/ms/axima-performance/index.html>.
- [25] G. J. McLachlan, *Discriminant analysis and statistical pattern recognition*. John Wiley and Sons, August 2004, vol. 544.
- [26] R. J. Tallarida and R. B. Murray, "Area under a curve: trapezoidal and simpson's rules," in *Manual of Pharmacologic Calculations*, 1987, pp. 77–81.
- [27] S. R. Safavian and D. Landgrebe, "A survey of decision tree classifier methodology," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 21, no. 3, pp. 660–674, May 1991.
- [28] L. Ceriani and P. Verme, "The origins of the gini index: extracts from variabilità e mutabilità (1912) by corrado gini," *The Journal of Economic Inequality*, vol. 10, no. 3, pp. 421–443, September 2012.
- [29] L. S. Shapley, "A value for n-person games," *Contributions to the Theory of Games* 2.28, no. 28, pp. 307–317, March 1953.