



# **TITLE : EXPLAINABLE AI**

**(ML PROJECT)**

**MANSURAH M - 22BAI1338**  
**GOWTHAMI S - 22BAI1457**



# RESEARCH PAPER

## Explainable Artificial Intelligence for Intrusion Detection System

Authors:

Shruti Patil, Vijayakumar

Varadarajan, Siddiqui Mohd Mazhar,

Submission received: 18 July 2022 /

Revised: 14 September 2022 /

Accepted: 20 September 2022 /

**Published:** 27 September 2022



electronics



Article

## Explainable Artificial Intelligence for Intrusion Detection System

Shruti Patil <sup>1,\*</sup>, Vijayakumar Varadarajan <sup>2,3,4,\*</sup>, Siddiqui Mohd Mazhar <sup>5</sup>, Abdulwodood Sahibzada <sup>5</sup>,  
Nihal Ahmed <sup>5</sup>, Onkar Sinha <sup>5</sup>, Satish Kumar <sup>1</sup>, Kailash Shaw <sup>5</sup> and Ketan Kotecha <sup>1</sup>

- <sup>1</sup> Symbiosis Centre for Applied Artificial Intelligence (SCAAI), Symbiosis Institute of Technology, Symbiosis International (Deemed University), Pune 412115, India
- <sup>2</sup> School of Computer Science and Engineering, The University of New South Wales, Sydney, NSW 1466, Australia
- <sup>3</sup> School of NUOVOS, Ajeenkya D Y Patil University, Pune 412105, India
- <sup>4</sup> Swiss School of Business and Management, 1213 Geneva, Switzerland
- <sup>5</sup> Department of Computer Science Engineering, Symbiosis Institute of Technology, Symbiosis International (Deemed University), Pune 412115, India
- \* Correspondence: shruti.patil@sitpune.edu.in (S.P.); vijayakumar.varadarajan@gmail.com (V.V.)



# OBJECTIVE

- The model incorporates the XAI algorithm LIME for better explainability and understanding of the black-box approach to intrusion detection system by observing in decision tree, random forest and support vector machine.
- The proposed IDS incorporates features from the CICIDS-2018 dataset .

# WORKFLOW OF THE PAPER:

## Workflow :

- 1.Data Loading
- 2.Data preprocessing
- 3.Feature Selection
- 4.ML model training
- 5.Applying LIME to the ML models (Decision tree, SVM, Random Forest classifier)
- 6.Generating user-understandable explanations

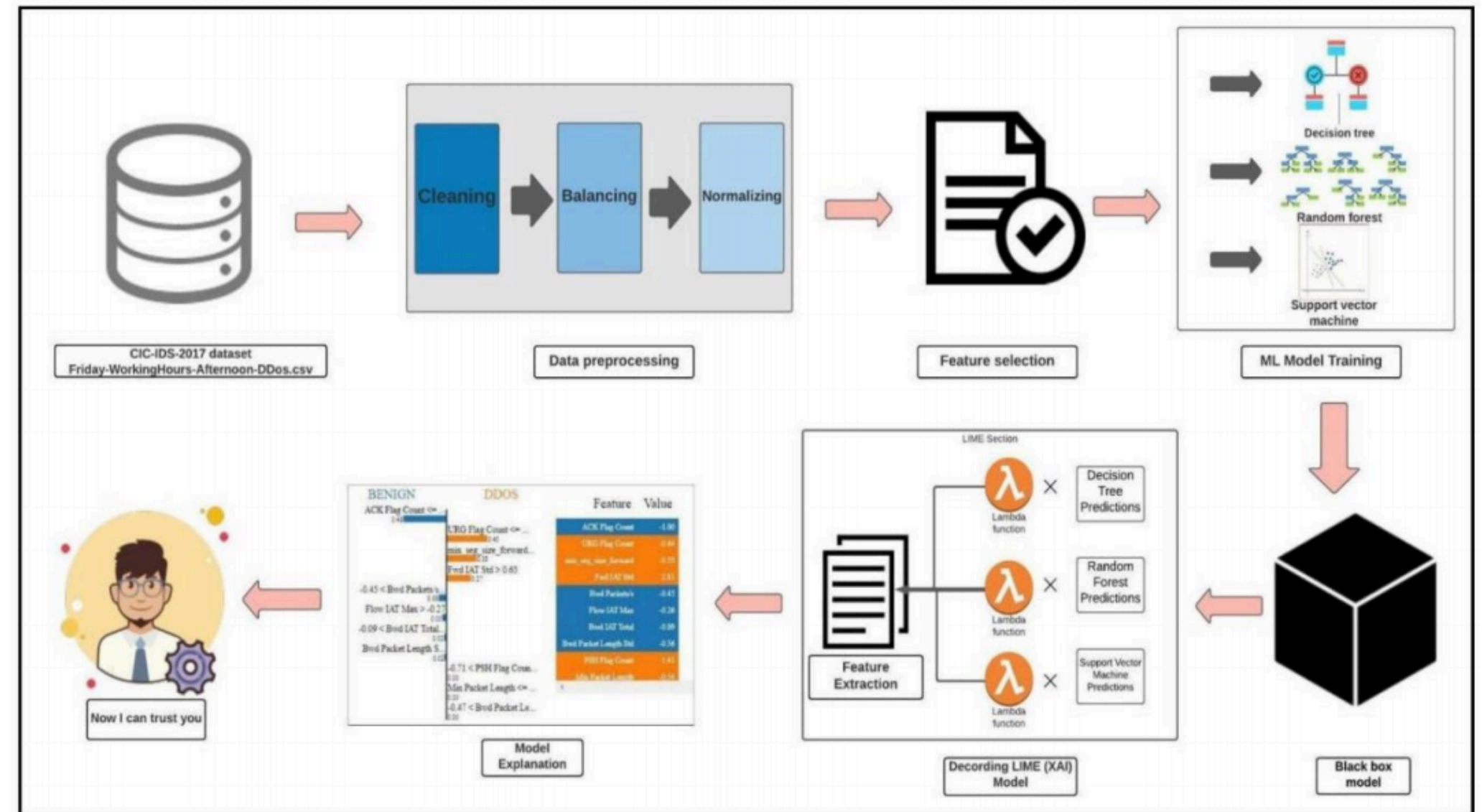
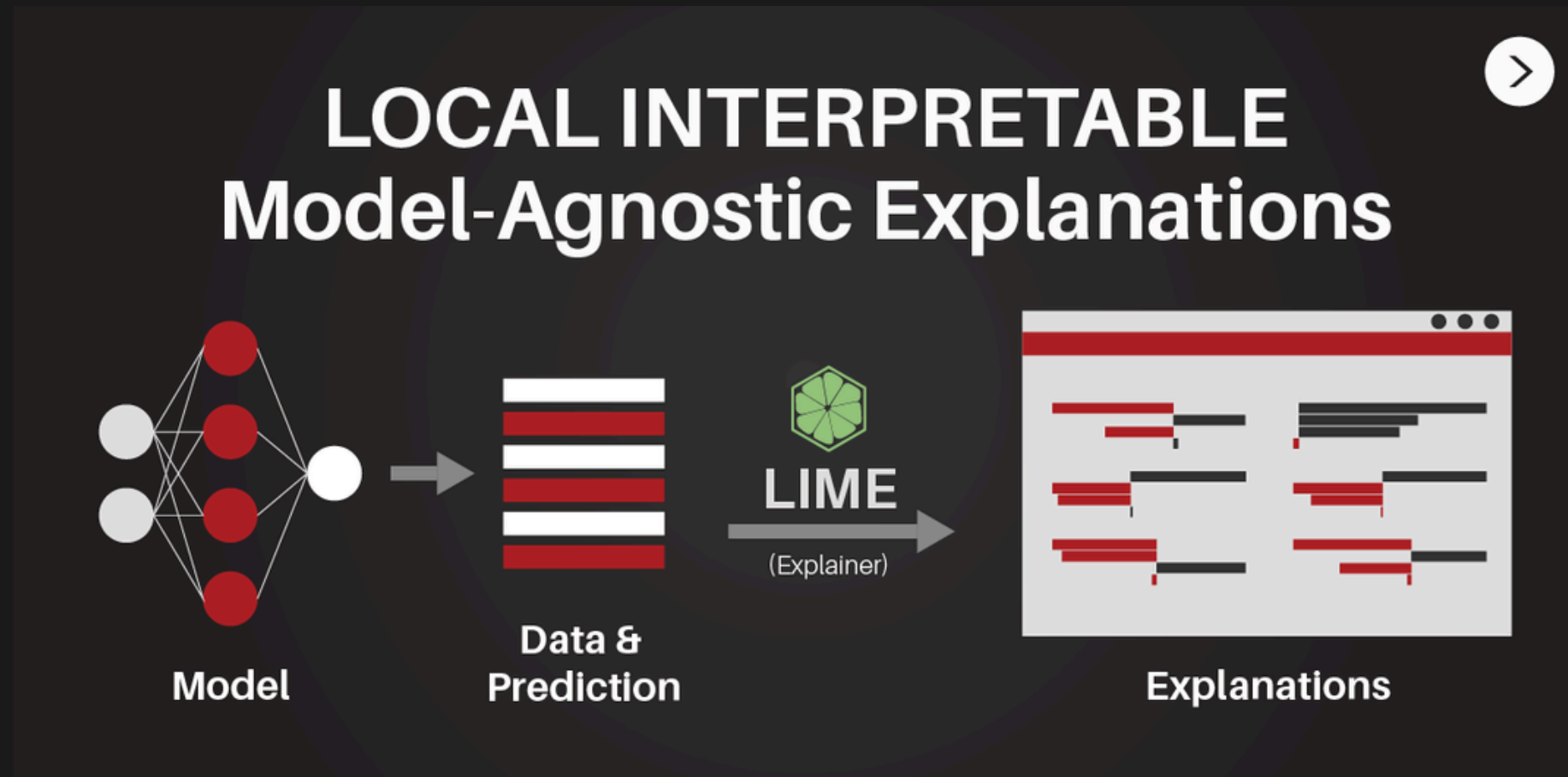


Figure 7. Explainable AI Framework.

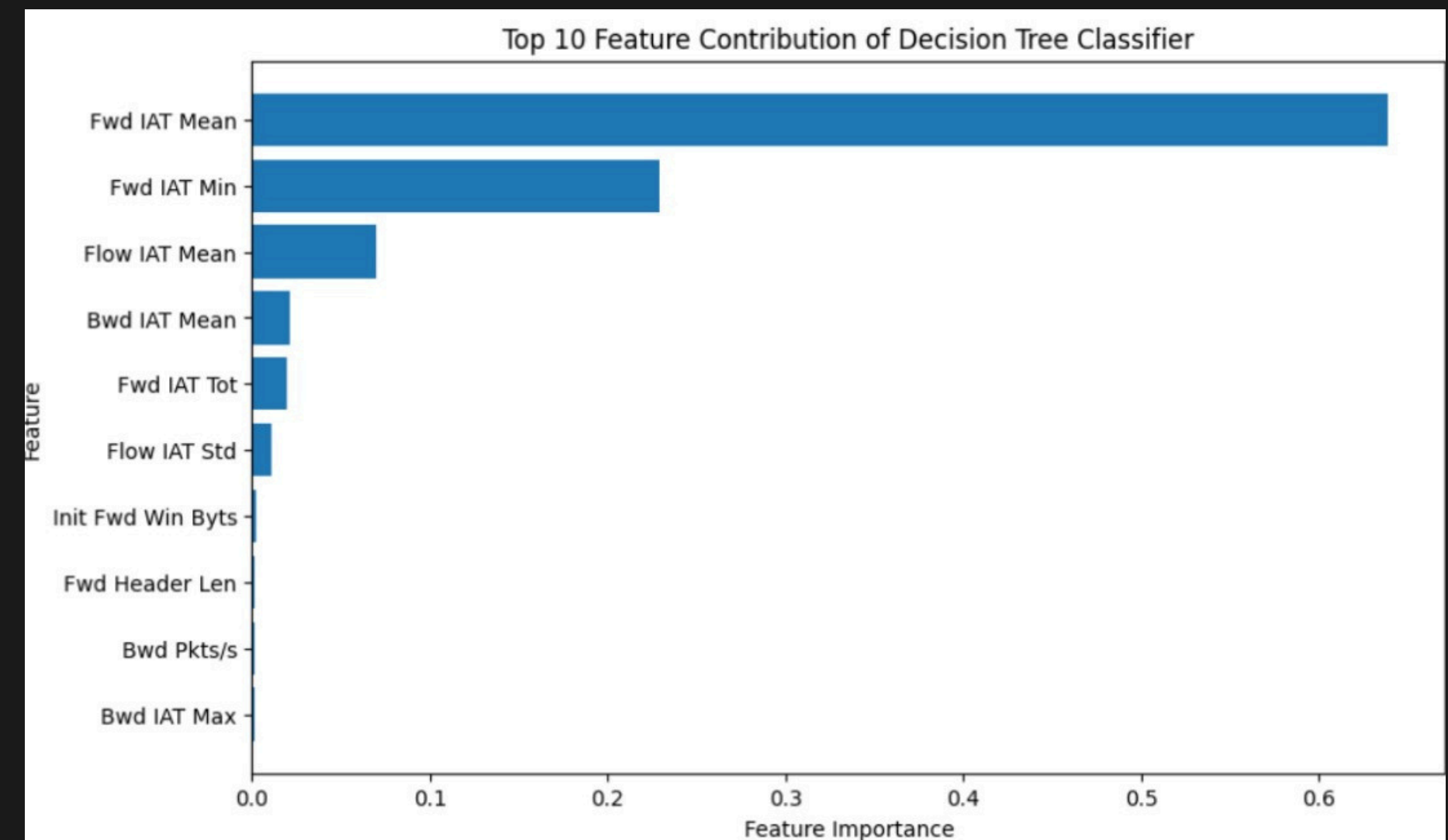
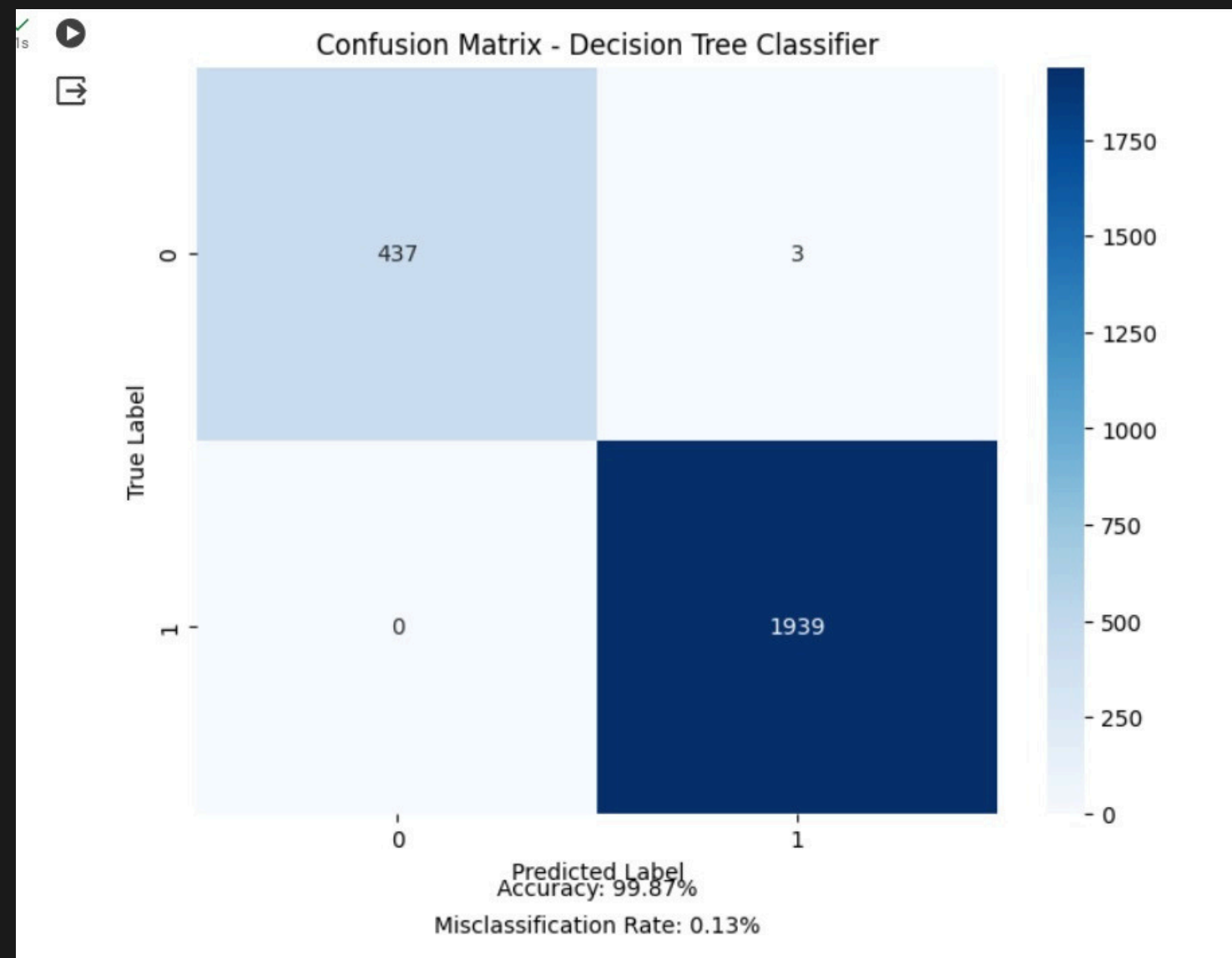


# LIME



- LIME, Local Interpretable Model-agnostic Explanations, is a technique that can be used to explain the predictions made by any black-box classifier.
- The proposed IDS incorporates features from the CICIDS-2018 dataset.

# DECISION TREE

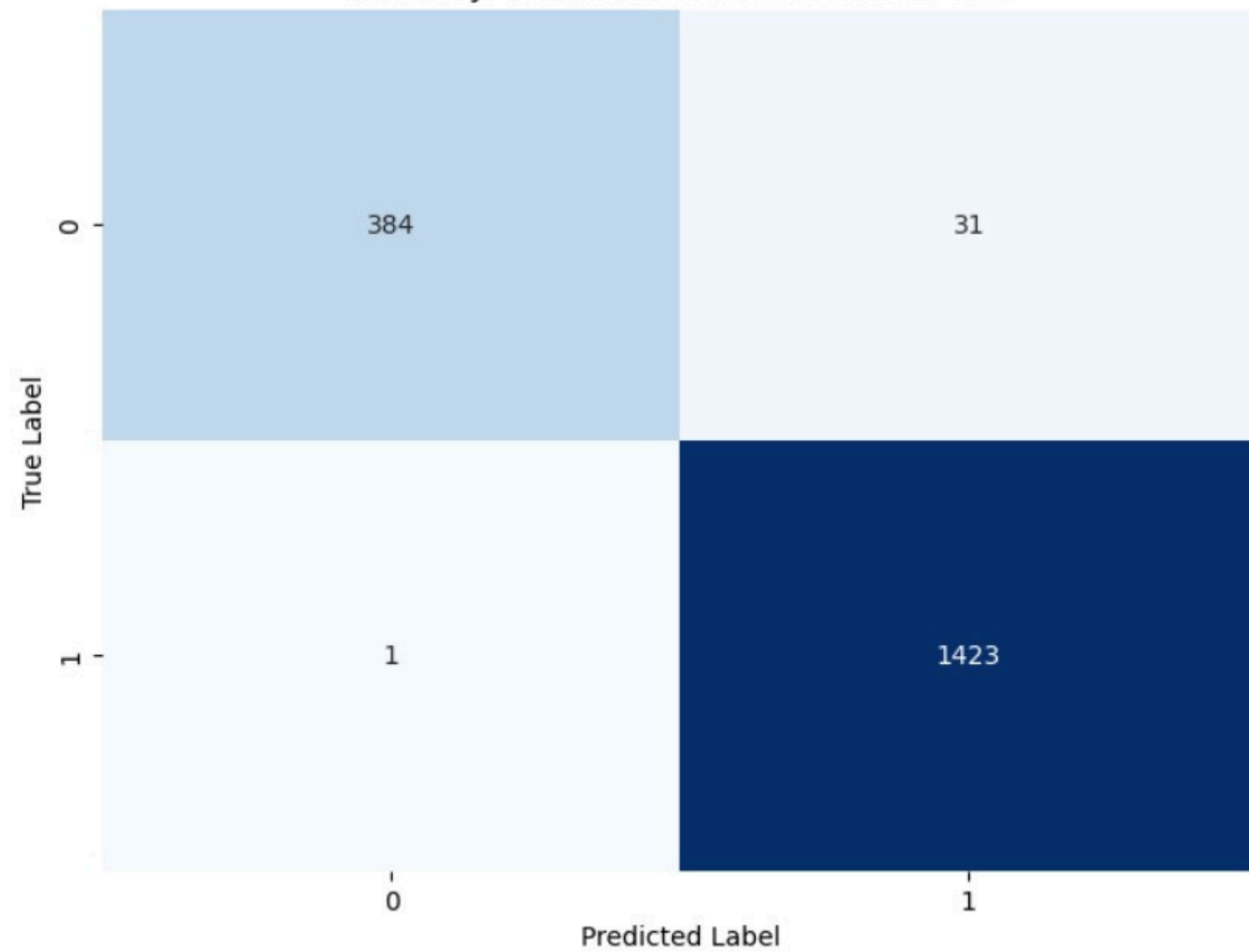


# GENERATING LIME EXPLANATIONS FOR DECISION TREE

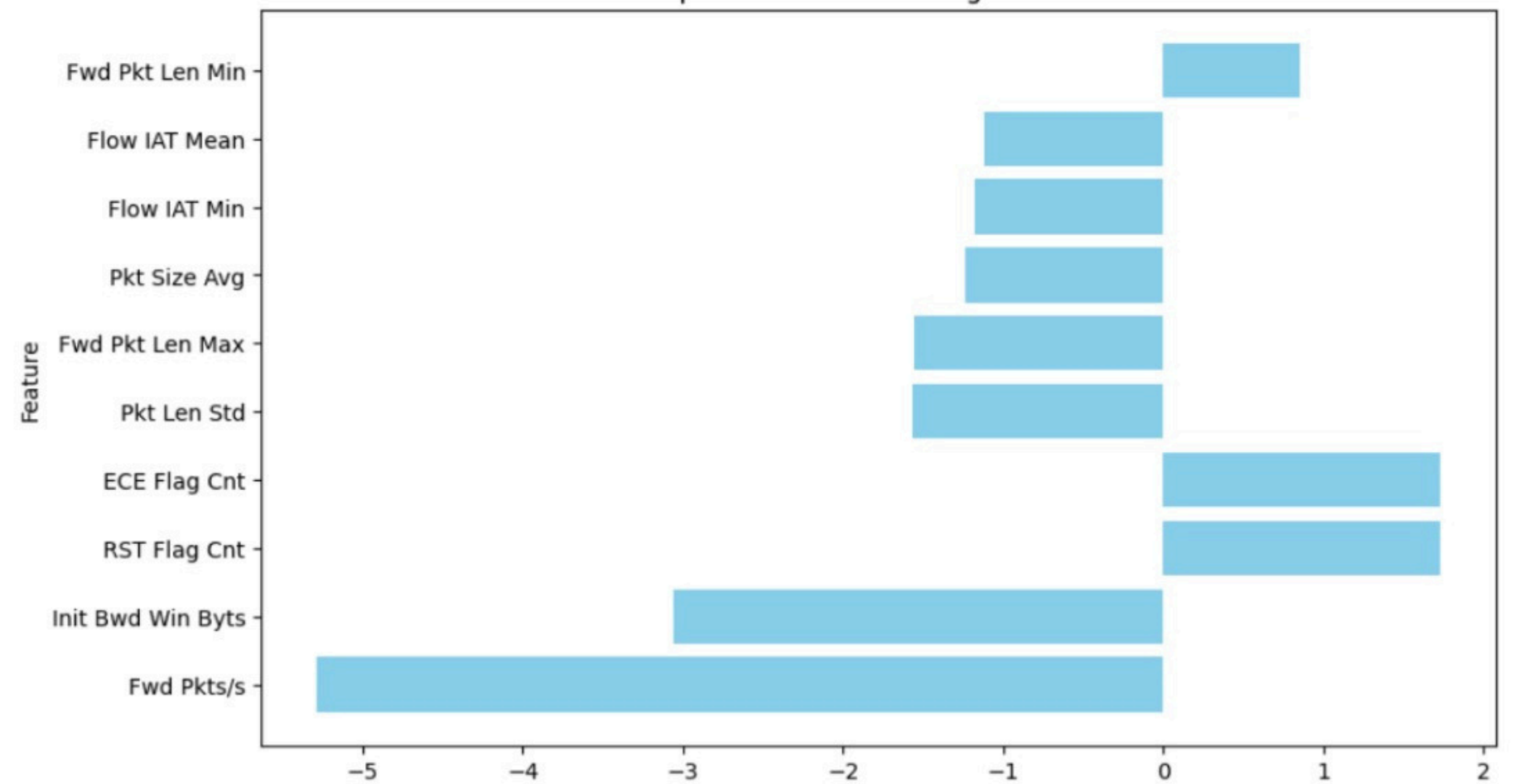


# SVM

Confusion Matrix  
Accuracy: 0.98, Misclassification Rate: 0.02



Top 10 Most Contributing Features





# SVM

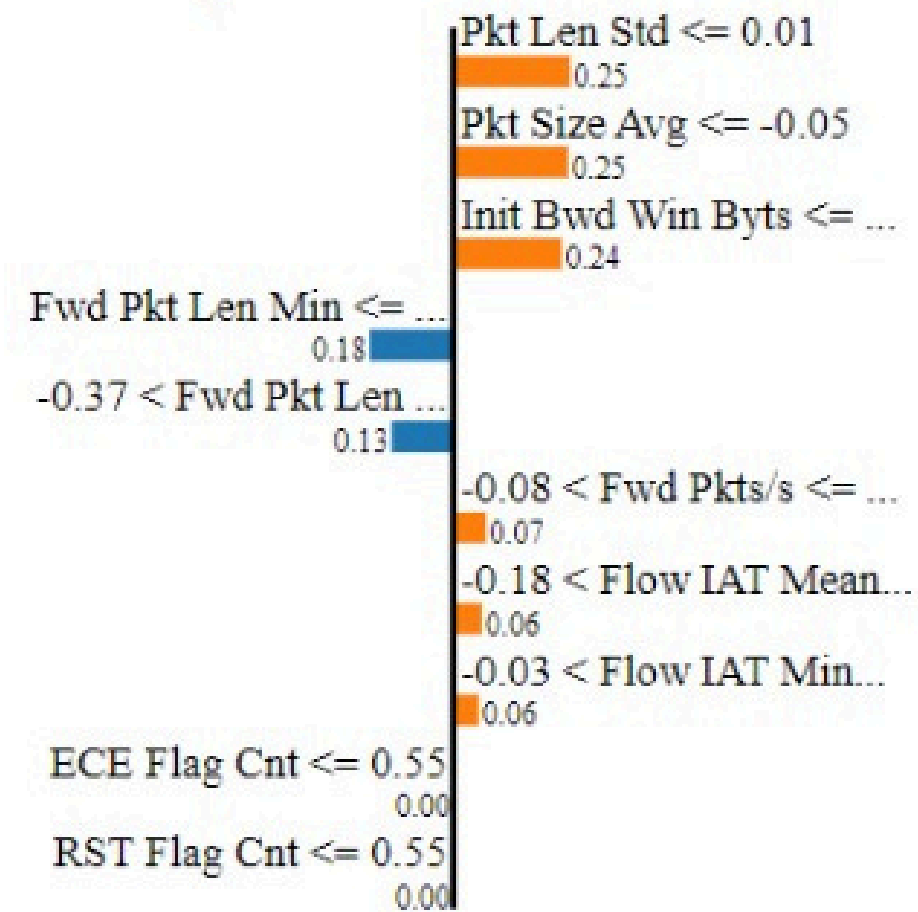


Prediction probabilities

Benign   
Bot 0.99

Benign

Bot

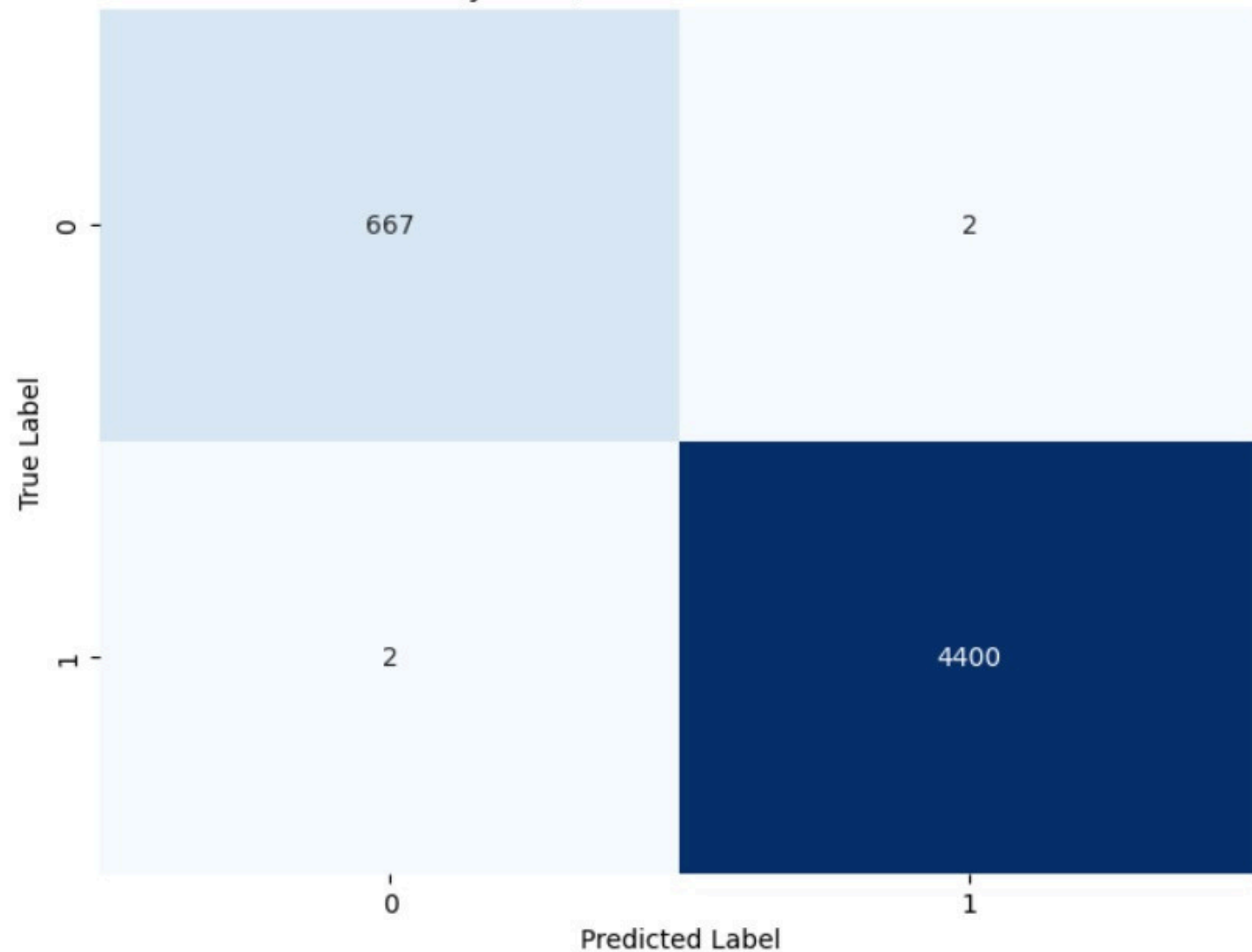


Feature Value

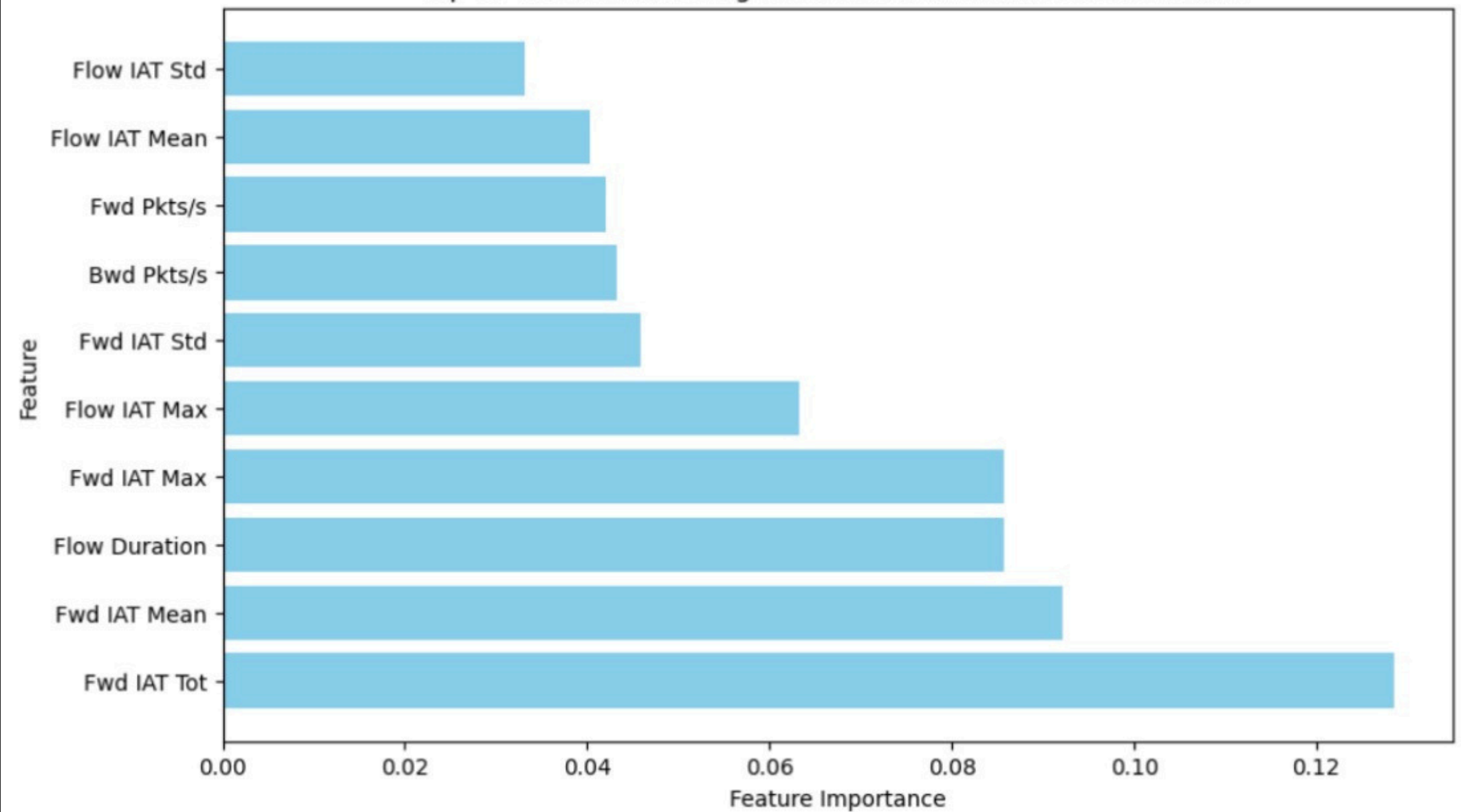
Pkt Len Std	0.01
Pkt Size Avg	-0.05
Init Bwd Win Byts	-0.17
Fwd Pkt Len Min	-0.10
Fwd Pkt Len Max	0.32
Fwd Pkts/s	-0.08
Flow IAT Mean	-0.18
Flow IAT Min	-0.03
ECE Flag Cnt	0.55
RST Flag Cnt	0.55

# RANDOM FOREST CLASSIFIER

Random Forest Confusion Matrix  
Accuracy: 1.00, Misclassification Rate: 0.00



Top 10 Most Contributing Features for Random Forest Classifier



# RANDOM FOREST CLASSIFIER

```
# Print LIME explanation  
rf_explanation_top10.show_in_notebook()
```

```
rf_classifier_top10.predict_proba,  
num_features=len(top10_features_rf))
```

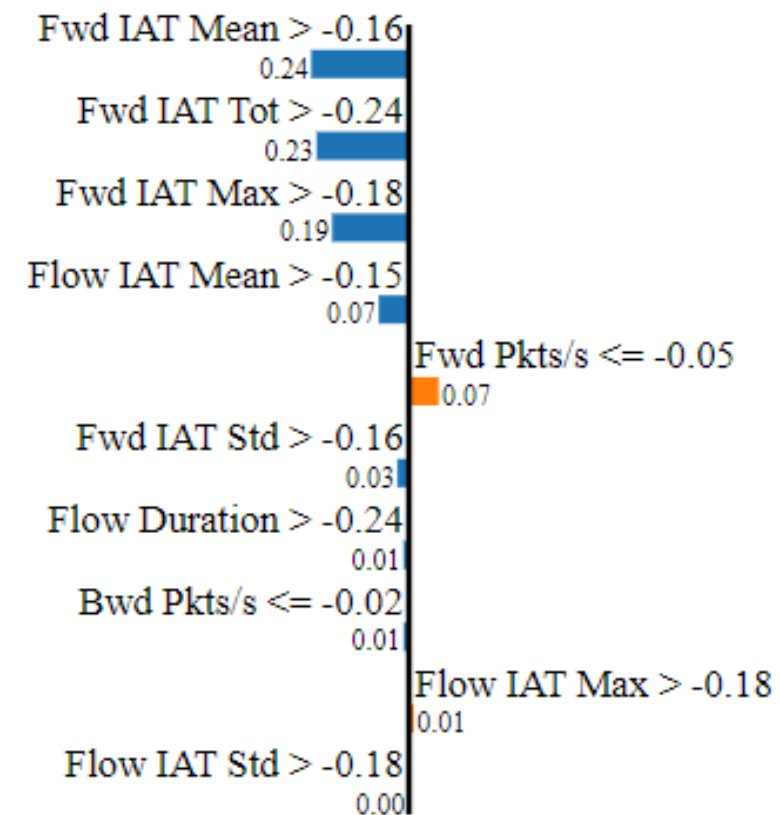


Prediction probabilities

Benign ☒ 1.00  
Bot ☐ 0.00

Benign

Bot



Feature Value

Fwd IAT Mean	-0.10
Fwd IAT Tot	-0.23
Fwd IAT Max	-0.14
Flow IAT Mean	-0.08
Fwd Pkts/s	-0.05
Fwd IAT Std	-0.08
Flow Duration	-0.23
Bwd Pkts/s	-0.03
Flow IAT Max	-0.14
Flow IAT Std	-0.10

**THANK YOU**

