CERIAS

The Center for Education and Research in Information Assurance and Security

XAI-ADS: An Explainable Artificial Intelligence Framework for Enhancing Anomaly Detection in Autonomous Driving Systems

Sazid Nazat¹, Lingxi Li¹, and Mustafa Abdallah¹

¹Indiana University-Purdue University Indianapolis, Indiana, USA snazat@iu.edu, II7@iupui.edu, mabdall@iu.edu

Introduction

• Autonomous vehicles (AVs) are vulnerable to many cyber-attacks, which might cause them to behave abnormally in vehicular ad hoc network (VANET). This disrupts the VANET.

• AI models are employed to classify anomalous AVs in the VANET, ensuring that it operates smoothly. XAI is then used to explain the decisions made by AI models in detecting anomalous AVs in a way that human operators (for example, safety drivers) can understand.

• To build **trust** between human operators and the AI models used for securing autonomous vehicles (AVs), XAI provides explanations to back up the AI predictions with evidence.

A) VeReMi datasetA) Verall Performance of AI modelsb) Sensor datasetCoveral performances for AI models with to factories for the VereMi datasetA) Veral Performances for AI models with to factories for the VereMi datasetA) Veral Performance for AI models with to factories for the VereMi datasetA) Veral Performance for AI models with to factories for the VereMi datasetA) Veral Performance for AI models with to factories for the VereMi datasetA) Veral Performance for AI models datasetA) Veral Performance for AI models datasetB) Sensor dataset	Results												
<section-header><section-header> a) Vereen eta Key insight: Table 5. 5. overall performances for AI models with to fatteres for the Vereen it ve</section-header></section-header>		Overall Performa	ance of AI models										
	A) Verence Section 2013TABLE 5: Overall performances for AI models with top 6 fatures for the VereMi dataset.<a href="https://doi.org/10.103/082008800880088008800880088008800880088</td> <td>Key insight: AdaBoost and RF classifiers give the best overall performance on VeReMi dataset.</td> <td><section-header><text><text></text></text></section-header></td> <td>Key insight: All Al models depicted high performance (except the SVM model) on Sensor dataset.</td>	Key insight: AdaBoost and RF classifiers give the best overall performance on VeReMi dataset.	<section-header><text><text></text></text></section-header>	Key insight: All Al models depicted high performance (except the SVM model) on Sensor dataset.									

DNN

Motivation

- The complicated feature interactions in AI models make it hard to grasp how and why they
 make decisions. For most AI models, this black-box issue makes it hard to explain their
 behavior, especially when mistakes arise.
- AI models have good prediction accuracy in various autonomous driving sectors, but their decision-making process must be interpreted and explained, especially in safety-critical domains.
- XAI can build trust and ensure accountability through explainable AI model decisions and insights. Also, it can enhance security and safety by understanding AI decision-making process and providing necessary information.

Our Contribution

- Propose an end-to-end XAI framework offering global and local explanations that provide human-understandable justifications for AI models' decisions in detecting anomalous AVs.
- Introduce **two novel XAI-based feature selection** approaches to identify and rank significant features contributing to anomaly classification of AVs and compare them with baseline feature selection methods, evaluated on **three independent classifiers** to ensure unbiased results.

 Benchmark the anomaly detection problem using six different AI models and two XAI methods on two autonomous driving datasets, analyzing the efficiency of the black-box models and the XAI techniques employed.

Framework

Beingn percentage	0010	0070	1010	0070	50%	0070	1010	0070	50%	0070	1010	0070		
Acc	0.74	0.74	0.75	0.76	0.77	0.77	0.78	0.78	0.50	0.49	0.54	0.51		
Prec	0.75	0.76	0.77	0.79	0.78	0.78	0.80	0.81	0.51	0.55	0.58	0.54		
Rec	0.76	0.77	0.82	0.84	0.80	0.79	0.83	0.85	0.64	0.68	0.73	0.67		
F-1	0.73	0.78	0.80	0.81	0.79	0.79	0.81	0.83	0.56	0.55	0.65	0.60		
		K	NN			SI	/ M		AdaBoost					
		111	111							nua	DUUSL			
Benign percentage	50%	60%	70%	80%	50%	60%	70%	80%	50%	60%	70%	80%		
Benign percentage Acc	50% 0.72	60% 0.72	70% 0.74	80% 0.75	50% 0.52	60% 0.55	70% 0.58	80% 0.61	50% 0.76	60% 0.76	70% 0.77	80% 0.78		
Benign percentage Acc Prec	50% 0.72 0.74	60% 0.72 0.74	70% 0.74 0.76	80% 0.75 0.76	50% 0.52 0.52	60% 0.55 0.55	70% 0.58 0.59	80% 0.61 0.61	50% 0.76 0.78	60% 0.76 0.79	70% 0.77 0.80	80% 0.78 0.82		
Benign percentage Acc Prec Rec	50% 0.72 0.74 0.68	60% 0.72 0.74 0.68	70% 0.74 0.76 0.71	80% 0.75 0.76 0.71	50% 0.52 0.52 0.55	60% 0.55 0.55 0.95	70% 0.58 0.59 0.99	80% 0.61 0.61 1	50% 0.76 0.78 0.73	60% 0.76 0.79 0.78	70% 0.77 0.80 0.80	80% 0.78 0.82 0.83		

RF

Key insight: AI models perform
better when the percentage of
anomalous data is reduced (less
number of anomalous samples)
except for the DNN model.

XAI Explanations vs Black Box AI Models



XAI Explanation Text: The total value of benign features is greater than that of the anomalous features. Therefore, the model predicted this sample correctly as benign.

(a) Correctly classified benign sample



XAI Explanation Text: Most of the features for this sample is anomalous. The model is confused because of the large negative value of the feature "spd_y". Note: "spd_y" is the second most influential feature in SHAP-AI feature selection method. The prediction probabilities are almost equal. That might be another reason for the model to get confused.

(a) Misclassified benign sample



TABLE 8: Performance of top-k features of SHAP-based and SHAP-AI feature selection methods on VeReMi dataset for the three independent classifiers (LR, LGBM, and XGBoost).

	Evaluation		LR	L	GBM	XG	Boost	
		SHAP-based	SHAP-AI based	SHAP-based	SHAP-AI based	SHAP-based	SHAP-AI based	
	Accuracy	0.50	0.50	0.67	0.67	0.67	0.67	
k-1	Precision	0.67	0.67	0.67	0.67	0.67	0.67	
K=1	Recall	0.49	0.49	0.99	0.99	0.99	0.99	
	F-1 score	0.56	0.56	0.80	0.80	0.80	0.80	
	Accuracy	0.54	0.53	0.67	0.67	0.68	0.68	
k-2	Precision	0.67	0.68	0.67	0.67	0.68	0.68	
K=2	Recall	0.61	0.56	0.99	0.99	0.98	0.97	
	F-1 score	0.64	0.61	0.80	0.80	0.80	0.80	
	Accuracy	0.56	0.57	0.68	0.67	0.69	0.69	
k-2	Precision	0.68	0.68	0.67	0.67	0.69	0.69	
к=э	Recall	0.65	0.67	0.99	0.99	0.97	0.97	
	F-1 score	0.66	0.68	0.80	0.80	0.81	0.81	
	Accuracy	0.54	0.54	0.68	0.68	0.70	0.70	
k-1	Precision	0.67	0.68	0.67	0.67	0.70	0.70	
K-4	Recall	0.60	0.61	1	1	0.97	0.97	
	F-1 score	0.63	0.64	0.80	0.80	0.81	0.81	
	Accuracy	0.54	0.54	0.67	0.67	0.69	0.69	
k-5	Precision	0.68	0.68	0.67	0.67	0.69	0.69	
к=Ј	Recall	0.60	0.61	1	1	0.97	0.97	
	F-1 score	0.63	0.64	0.80	0.80	0.81	0.81	
	Accuracy	0.54	0.54	0.67	0.67	0.69	0.69	
k-6	Precision	0.67	0.67	0.67	0.67	0.69	0.69	
K=0	Recall	0.60	0.60	1	1	0.97	0.97	
	F-1 score	0.64	0.64	0.80	0.80	0.81	0.81	

Key insight: To rigorously evaluate the proposed feature selection approaches without bias, we used three independent classifiers, logistic regression (LR), light gradient boosting machine (LGBM), and extreme gradient boosting (XGBoost), which were not involved in feature ranking and performed well.



Evaluation



Sensor Name Formality 1bit – 10bit Checks every message if it is maintaining correct formality 0/1 Checks if the message reached the destined location Location 1Hz – 10Hz Checks the interval time of messages Frequency Checking if the AV is in the speed limit (Highway) 50mph - 90mph Checks if several messages adheres to defined specification 0/1 Correlation Lane Alignment | 1-3 Checks if the AV is in the correct lane Checks if the AV maintains the headway time range Headway Time 0.3s - 0.95s Checks for the correct order of communication messages Protoco 1 - 10000Plausibility 50% - 200% Checks if the data is plausible (relative difference between sizes of two consecutive payloads) Checks if all the parts of the AV are delivering consistent information about an inciden Consistency

LIVIE ON SECONSOL CLASSED abilities 1.00 Benign Anomalous Consistency > 0.87 Plausibility <= ... 0.03 Correlation > 0.58 0.02 Correlation > 0.58 0.02 Speed > -0.03 0.02 Formality > 0.76 Potocol > 0.85 0.02 Formality > 0.76 0.03 Correlation > 0.58 0.02 Formality > 0.76 0.02 Formality > 0.76 0.03 Correlation > 0.58 Correl

(a) Example of a local explanation for correctly predicting an anomalous AV from the Sensor dataset.

Comparison with Baseline Feature Selection Models

TABLE 13: Results per class in every feature selection method on VeReMi dataset under k = 4 for feature selection methods.

Model		SHAP					SHAP-AI			K-best			PCA				ET				IG				
	Class	Acc	Prec	Rec	F-1	Acc	Prec	Rec	F-1	Acc	Prec	Rec	F-1	Acc	Prec	Rec	F-1	Acc	Prec	Rec	F-1	Acc	Prec	Rec	F-1
DT	Benign (0)	0.76	0.78	0.90	0.83	0.79	0.82	0.87	0.84	0.78	0.82	0.87	0.84	0.78	0.82	0.87	0.84	0.78	0.82	0.87	0.84	0.78	0.81	0.87	0.84
	Anomalous (1)	0.70	0.70	0.49	0.50	0.75	0.70	0.61	0.66	0.70	0.70	0.61	0.65	0.70	0.70	0.61	0.65	0.70	0.70	0.61	0.65	0.70	0.70	0.60	0.65
RF	Benign (0)	0.80	0.83	0.88	0.86	0.80	0.83	0.88	0.85	0.80	0.83	0.88	0.86	0.80	0.83	0.88	0.85	0.80	0.83	0.88	0.86	0.80	0.83	0.88	0.85
	Anomalous (1)	0.00	0.73	0.64	0.68	0.00	0.72	0.64	0.68	0.00	0.73	0.64	0.68	0.00	0.73	0.64	0.68	0.00	0.73	0.64	0.68	0.60	0.73	0.63	0.68
DNN	Benign (0)	0.66	0.67	0.96	0.79	0.41	0.67	0.24	0.36	0.67	0.67	1	0.80	0.54	0.67	0.64	0.65	0.65	0.67	0.96	0.79	0.64	0.67	0.92	0.77
	Anomalous (1)	0.00	0.33	0.04	0.07	0.41	0.33	0.75	0.46	0.07	0.33	0.27	0.29	0.04	0.33	0.36	0.34	0.00	0.33	0.04	0.08	0.04	0.30	0.07	0.11
KNN	Benign (0)	0.78	0.82	0.85	0.84	0.79	0.82	0.85	0.84	0.79	0.82	0.85	0.84	0.79	0.82	0.85	0.84	0.79	0.82	0.85	0.84	0.78	0.82	0.85	0.83
	Anomalous (1)	0.70	0.68	0.63	0.65	0.78	0.68	0.63	0.65	0.78	0.68	0.64	0.66	0.76	0.68	0.63	0.65	0.78	0.68	0.63	0.65	0.78	0.67	0.63	0.65
SVM	Benign (0)	0.55	0.54	0.69	0.61	0.55	0.54	0.66	0.60	0.67	0.69	0.93	0.79	0.65	0.68	0.92	0.78	0.66	0.69	0.90	0.78	0.67	0.69	0.92	0.79
	Anomalous (1)	0.55	0.57	0.41	0.48	0.55	0.57	0.42	0.48	0.07	0.50	0.14	0.22	0.05	0.50	0.14	0.22	0.00	0.50	0.15	0.23	0.07	0.49	0.16	0.24
AdaBoost	Benign (0)	0.70	0.84	0.85	0.85	0.80	0.83	0.86	0.85	0.79	0.83	0.86	0.85	0.70	0.84	0.86	0.85	0.78	0.84	0.85	0.84	0.79	0.84	0.85	0.84
	Anomalous (1)	0.75	0.70	0.65	0.67	0.00	0.70	0.66	0.68	0.19	0.70	0.65	0.68	0.75	0.70	0.66	0.68	0.76	0.70	0.66	0.62	0.19	0.69	0.66	0.68

Key insight: Compared to six baseline methods across AI models, the novel feature selection methods excel in DT and AdaBoost models. They also improve precision, recall, and F1-score for anomaly detection in most of the models.

<section-header>a) VeReMidataset <u>× toroto viz viz viz vizi vizi vizi vizi vizi v</u></section-header>		Efficiency of XAI models													
$\frac{\text{XAI model}}{\text{SHAP}} = \frac{\text{VR}}{500} \frac{\text{VR}}{0.72} \frac{\text{RF}}{0.49} \frac{\text{DNN}}{0.13} \frac{\text{KNN}}{2.47} \frac{\text{SVM}}{19.09} \frac{\text{ADA}}{1.13} \\ \frac{500}{1.12} \frac{2.47}{0.49} \frac{91.03}{0.13} \frac{2.47}{0.49} \frac{91.09}{1.126} \frac{4.15}{4.77.23} \\ \frac{11}{10k} \frac{219.66}{535.31} \frac{535.31}{6.98} \frac{6.98}{929.98} \frac{929.98}{6711.26} \frac{477.23}{477.23} \\ \frac{500}{50k} \frac{8054.65}{13256.81} \text{NA} \frac{22128.2}{271275} \frac{271275}{8828.6} \\ \frac{500}{10k} \frac{3256.81}{10k} \frac{\text{NA}}{33.11} \frac{5.66}{5.66} \frac{51.6}{24.71} \frac{4.77}{0.65} \\ \frac{11}{1k} \frac{0.43}{0.48} \frac{500}{0.17} \frac{0.17}{0.39} \frac{3}{3} \frac{5.98}{0.48} \frac{0.48}{0.31} \\ \frac{11k}{1k} \frac{0.68}{0.167} \frac{1.17}{0.52} \frac{13.97}{1.397} \frac{0.74}{0.74} \frac{11.32}{1.32} \\ \end{array} \right)$	a) V	'eRe	₽Mi	da	tas	et			Key insight: LIME has better efficiency	b) S	Senso	r dat	aset		Key insight: LIME
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	SHAP	# Samples 500 1k 10k 50k	0.72 2.97 219.66 8054.65	1.84 6.29 535.31 13256.81	0.13 0.15 6.98 NA	KNN 2.47 9.06 929.98 22128.2	8VM 19.09 93.66 6711.26 271275	ADA 1.13 4.15 477.23 8828.6	than SHAP for	XAI model SHAP	# Samples D1 500 10.2 1k 33.2 10k 3826	RF 22 20.09 56 121.20 .31 8291	DNN KNN 0.08 329.04 0.15 866.01 8.20 97257.36 9	AdaBoost 20.43 10.47 51.29 50.28 513.15 3391.36	for Sensor dataset
	LIME	500 1k 10k 50k	0.12 0.43 27.24 27.23	2.48 3.11 43.1 131.88	4.65 5.66 95.90 501.67	4.94 5.16 89 287.75	25.53 24.71 379.38 NA	1.04 0.65 18.88 56.41	majority of the Al	LIME	500 0.1 1k 0.6 10k 10.2	7 0.39 8 1.17 27 19.39	3 5.98 6.25 13.97 66.56 75.46	0.48 0.31 0.74 1.15 17.64 11.32	compared to SHAP.

Acknowledgement

This work was supported by Lilly Endowment (AnalytixIN); Enhanced Mentoring Program with Opportunities for Ways to Excel in Research (EMPOWER), and the 1st Year Research Immersion Program (1RIP) Grants from the Office of the Vice Chancellor for Research at Indiana University – Purdue University Indianapolis.



