

False Data Injection Detection in Power System Based on LOSSA-AdaBoostDT

Lei Xi, *Member, IEEE*, Xilong Tian, Miao He, and Chen Cheng

Abstract—The attack of false data injection can contaminate the measurements acquired from the supervisory control and data acquisition (SCADA) system, which can seriously endanger the safety and stability of power system operations. The conventional machine learning attack detection methods use a single strong classifier and are difficult to solve the problem of overfitting, making them lack of generalization ability. On the other hand, most existing dimension reduction approaches based on feature extraction can change the original physical meanings of measurements. Here, a novel method is proposed based on feature selection and ensemble learning to solve the above problems. Squirrel search algorithm combines Latin hypercube sampling and opposition-based learning to form an improved algorithm with strong global search ability for feature selection. This avoids the problem of feature extraction changing the original physical meanings of measurements. Besides, the classifier based on adaptive boosting decision tree ensemble learning algorithm with stronger generalization ability is used to distinguish the false data injection. Simulation results using the IEEE 14-bus and IEEE 57-bus test systems verify the proposed method with higher performance of detection compared with other widely adopted methods.

Index Terms—False data injection, squirrel search algorithm, adaptive boosting decision tree, SCADA.

I. INTRODUCTION

With highly integrated computer, communication and control technologies, traditional power systems have developed into a new type of cyber-physical power systems (CPPS) [1]–[4]. In the CPPS, measurements obtained from the supervisory control and data acquisition (SCADA) and the phasor measurement units (PMU) of wide area monitoring systems are used to acquire the system operation states [5], [6], as reference for stability analysis, safety constraint scheduling and other performance analysis. The exactitude of the states estimation is largely dependent on the reliability of the measurements. Therefore, certain damaged data detection mechanism [7] should be adopted to identify the bad data in the measurements. However, due to the attacks on the SCADA and PMU measurements, the false data initiated by the false data injection attack (FDIA) [8], [9] can bypass the data scrutiny system so that the measurements are contaminated in the course of data collection, transmission and processing, seriously jeopardizing the safety of the power system [10], [11].

Many investigations have been carried out for FDIA detection in CPPS, which can be divided into three categories, i.e., improved state estimation [12]–[14], trajectory analysis [15]–[17] and machine learning based methods [18]–[20]. Both the improved state estimation and trajectory analysis methods are based on physical models, which depend on the specific attack model and system topology. The detection methods based on physical models are difficult to expand to other applications when the detection scenario changes. FDIA detection can be treated as an anomaly detection problem [21] via classifier utilizing historical measurements in machine learning methods. Strong classifier is a method with generalization ability for unknown data [22], which can classify the measurements into normal states and abnormal states more quickly and accurately with no knowledge of specific attack model and system topology, showing high scalability and great advantages in the FDIA detection. Reference [18] uses convolutional neural network to detect FDIA, and combines gradient descent strategy and classification cross-entropy loss to improve the convolutional neural network, which can identify a variety of abnormal events. Reference [19] proposes a detection scheme based on a hybrid chimp optimized extreme learning machine to detect FDIA,

Received: July 9, 2024

Accepted: November 5, 2024

Published Online: May 1, 2025

Lei Xi (corresponding author) is with the College of Electrical Engineering & New Energy, and Hubei Provincial Key Laboratory for Operation and Control of Cascaded Hydropower Station, China Three Gorges University, Yichang 443002, China (e-mail: xilei2014@163.com).

Xilong Tian is with Qianjiang Power Supply Company, State Grid Hubei Electric Power Co., Ltd., Qianjiang 433100, China (e-mail: tianxilong1024@163.com).

Miao He is with Jingzhou Power Supply Company, State Grid Hubei Electric Power Co., Ltd., Jingzhou 434000, China (e-mail: he_miao98@163.com).

Chen Cheng is with Southern Power Grid Co., Ltd., Foshan Supply Bureau, Leliu Power Supply Station, Foshan 528000, China (e-mail: 842125376@qq.com).

DOI: 10.23919/PCMP.2024.000129

while in [20], the graph edge-conditioned convolutional networks is used to improve the detection framework to realize the FDIA detection and other network attacks.

However, the above mentioned machine learning based FDIA detection methods are constructed with single strong classifier [23], which is difficult to tackle the overfitting during model training, as a bottleneck for generalization capability improvement. Compared with the single strong classifier, ensemble strong classifier is constructed by multiple weak classifiers, i.e., ensemble learning [24], to provide regularization so as to alleviate overfitting. Further, a bootstrap aggregation scheme has been used to form an ensemble strong classifier to promote the classifier performance [25]. Reference [26] introduces an ensemble strong classifier based on extremely randomized ensemble learning which can detect the exact position of FDIA without statistical knowledge assumption. Reference [27] proposes to use extreme gradient boosting ensemble learning to detect non-stationary and nonlinear attacks. However, the current ensemble learning can only reduce one of the bias or variance of the ensemble strong classifier at one time to limit the diminution of the generalization error.

Adaptive boosting (AdaBoost) [28] is a kind of ensemble learning algorithm which can reduce the variance and bias of the classifiers simultaneously with better classification performance. In [29], AdaBoost based on support vector machine is applied to detect the chatter vibration with a high accuracy, while reference [30] achieves better accuracy of high-voltage circuit breaker fault diagnosis under small samples and complex working conditions by utilizing AdaBoost ensemble learning. Further, AdaBoost based on back propagation neural networks is proposed to identify rotor fault status and fault severity [31]. However, the improvement of generalization performance of the ensemble learning is to obtain the diversity among the weak classifiers, and thus, it is difficult to increase the diversity degree in the above methods. Decision tree is a classifier with strong interpretability, so a set of diversified weak classifiers can be obtained by adjusting the tree depth and width to acquire an ensemble with stronger generalization. Further, reference [32] uses the decision tree as the weak classifier in adaptive boosting ensemble learning to develop the adaptive boosting decision tree (AdaBoostDT), and achieve better computational efficiency in electricity theft detection. Therefore, this paper attempts to explore a classifier based on AdaBoostDT algorithm to detect the FDIA in CPPS with higher accuracy and fast speed.

It is worth noting that CPPS is a large-scale complex network with high-dimensional generated measurements [33]. Using high dimensional data to train the classifier will lead to a high complexity of the model, while the redundant features will weaken the detection performance of the model. Therefore, the dimension reduction [34] technology should be adopted in the original measurements to improve the efficiency of the classifier.

At present, there are limited studies on the data dimension reduction in FDIA detection. In [35], the autoencoder is integrated into the generative adversarial network and feature extraction dimension reduction is performed on the measurements. Reference [36] uses the principal component analysis to perform feature extraction and develops visual clustering method to detect FDIA. However, these dimension reduction methods are based on the feature extraction strategy which can damage the original representation and physical meaning of the data features.

Reference [37] introduces a feature selection strategy based on genetic algorithm to eliminate the redundant features so as to enhance the detection ability of the classifier. In [38], feature selection is performed to select the optimal feature subset from the original features, in which the original physical features of the units do not change, i.e., the measurements retain their original meanings. However, genetic algorithm such evolutionary heuristic algorithms have defects, i.e., slow search speed, heavy computation load with high dimensions data and low convergence.

Reference [39] designs the squirrel search algorithm (SSA) with fast convergence and strong adaptability to high dimensional data. However, like other swarm intelligence optimization algorithms, SSA is prone to be caught in local optimum [40]. In [41], the Latin hypercube sampling (LHS) is introduced to initialize the population, which can diversify the variety and ergodicity of the population, while reference [42] proposes the opposition-based learning (OBL) to improve the global search capability of the algorithm. To further optimize SSA for the feature subset selection, this paper integrates LHS and OBL into SSA to obtain the Latin hypercube sampling and opposition-based learning SSA (LOSSA) which can improve the global search capability so as to avoid the local optimum.

A novel FDIA detection method in CPPS based on LOSSA-AdaBoostDT algorithm is proposed in this paper. The main contributions of the work are summarized as follows.

- 1) In order to avoid the alternation of the original physical meaning of the measurements during the feature extraction, an integrated LOSSA scheme is proposed to combine the LHS and OBL strategies in SSA for feature selection of the measurement without trapped in local optimum.

- 2) A classifier based on AdaBoostDT algorithm with stronger generalization ability is designed to detect FDIA more quickly and accurately.

- 3) The comprehensive simulations of IEEE 14-bus and IEEE 57-bus test systems are used in the experiment to verify the effectiveness of the proposed FDIA detection method.

The remainder of the paper is arranged as follows. Section II explains the mechanism of FDIA, while Section III presents the proposed LOSSA-AdaBoostDT. In

Section IV, the experiments and numerical results are analyzed. Finally, Section V concludes the paper.

II. FALSE DATA INJECTION ATTACK MODEL

The state estimation can map between the real-time measurements and the system state variables. In order to simplify the state estimation, linearized direct current power flow model is often used to approximate alternating current model in engineering. The mapping between the measurements and the system states is described as:

$$\mathbf{m} = \mathbf{H}\mathbf{s} + \mathbf{e} \quad (1)$$

where \mathbf{m} is the system measurement vector; \mathbf{H} is the Jacobian matrix of the power system; \mathbf{s} is the state vector which includes the voltage magnitude and phase angle; and \mathbf{e} is the measurement error. It is noted that the measurement error includes measurement noise, communication noise and other noise.

The attack model of the FDIA can be described as:

$$\mathbf{m}_f = \mathbf{m} + \mathbf{f} = \mathbf{H}\mathbf{s} + \mathbf{f} + \mathbf{e} \quad (2)$$

where \mathbf{m}_f represents the measurement vector after the attack; and \mathbf{f} is the false data vector injected by the attacker. FDIA can bypass the general false data detection mechanism, and the false data vector would meet the constraint:

$$\mathbf{f} = \mathbf{H}\mathbf{b} \quad (3)$$

where \mathbf{b} is the bias of the state variable produced by the attack vector.

After the fake data are injected, the attacked measurements can be rewritten as:

$$\mathbf{m}_f = \mathbf{H}\mathbf{x} + \mathbf{H}\mathbf{b} + \mathbf{e} = \mathbf{H}(\mathbf{x} + \mathbf{b}) + \mathbf{e} = \mathbf{H}\mathbf{x}_f + \mathbf{e} \quad (4)$$

where \mathbf{x}_f is the new system state variables containing the false data after being attacked.

After the FDIA attack, the residual of the false data detection \mathbf{r}_f of the state estimation is written as:

$$\mathbf{r}_f = \mathbf{m}_f - \mathbf{H}\mathbf{x}_f = \mathbf{m} + \mathbf{f} - \mathbf{H}(\mathbf{x} + \mathbf{b}) + \mathbf{e} = \mathbf{H}\mathbf{x}_f + \mathbf{e} \quad (5)$$

Thus, FDIA can steer by the conventional residual-based false data detection and jeopardize the power system state estimation [43]. The FDIA structure in the CPPS is presented in Fig. 1.

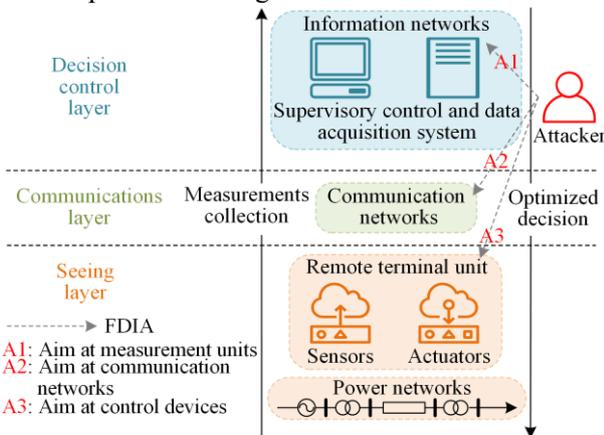


Fig. 1. FDIA structure.

It can be seen that the attacker can inject false data in the control equipment, communication networks and measurement units, to temper the measurement values.

III. LOSSA-ADABOOSTDT FALSE DATA DETECTION

A detection algorithm based on LOSSA-AdaBoostDT is proposed to detect FDIA in the CPPS. The LOSSA algorithm can perform the feature selection on the dataset generated by the attack model to reduce the dimension of the original measurement features as a data pre-processing algorithm, and AdaBoostDT is used to construct a classifier for FDIA detection. The main process of the FDIA detection algorithm is illustrated as in Fig. 2. The proposed detection model does not depend on the specific attack propagation mode and system model, and learns the outlier characteristics caused by FDIA in the electrical force measurements through training, so as to achieve the attack detection. It is suitable for FDIA which is aimed at CPPS state estimation in various scenarios. The detection model is deployed before advanced applications such as online power flow, security analysis, and economic scheduling, and when FDIA is detected, the state estimation results are replaced or restored to ensure that the energy management system makes correct decision on control instructions.

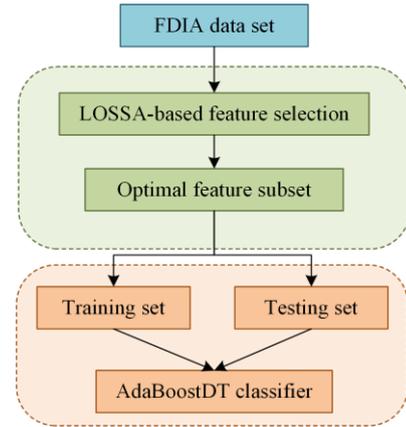


Fig. 2. The procedure of LOSSA-AdaBoostDT detection method.

A. Squirrel Search Algorithm

Compared with the dimension reduction method of feature extraction, feature selection is used to choose the optimal feature subset from all the original features and interpretably preserve the physical meaning of the data. It is a nondeterministic polynomial time combinatorial optimization problem, so SSA is applied for feature selection due to its simple principle, fast convergence and efficient processing of large-scale high-dimensional data. To simulate the forage process and special gliding mode of the squirrels, four basic hypotheses are made:

- 1) There are n squirrels in the forest with n trees.
- 2) The forest contains one hickory tree and N_s ($1 < N_s < n$) oak trees, and the remainder are the ordinary trees.

3) Among these three kinds of trees, the hickory tree has the best food, the oak trees have the common food and the ordinary trees have no food.

4) Each squirrel individually searches for food.

Under the above assumptions, the whole process of SSA is as follows.

Step 1: Population initialization

The squirrels are initialized as:

$$\mathbf{F}_i = \mathbf{F}_L + \text{rand}(1, D) \times (\mathbf{F}_U - \mathbf{F}_L) \quad (6)$$

where \mathbf{F}_i is the position of the i th ($i=1,2,\dots,n$) squirrel; \mathbf{F}_U and \mathbf{F}_L are the upper and lower bounds of the search space respectively; $\text{rand}(\cdot)$ is a random number in the range $[0,1]$; and D is the dimension of the individual.

Step 2: Individual classification

After the population initialization, the fitness (f) of the position of each individual is calculated and sorted in ascending sequence, i.e.:

$$\mathbf{F}_{\text{index}} = \text{sort}(f) \quad (7)$$

where $\mathbf{F}_{\text{index}}$ is all \mathbf{F} sorted in ascending by f . Then, the squirrels are divided into three categories based in their fitness values:

$$\mathbf{F}_h = \mathbf{F}_{\text{index}}(1) \quad (8)$$

$$\mathbf{F}_a(1:3) = \mathbf{F}_{\text{index}}(2:4) \quad (9)$$

$$\mathbf{F}_n(1:n-4) = \mathbf{F}_{\text{index}}(5:n) \quad (10)$$

where \mathbf{F}_h is the squirrel number in the hickory trees; \mathbf{F}_a is the squirrel number in the oak trees; and \mathbf{F}_n is the squirrel number in the normal trees.

Step 3: Location update

The squirrels use a special guidance to search trees and update their positions, while the presence of the predators can affect such behaviors. Three situations can occur during the search process, which are shown as follows.

1) Squirrels move from the oak trees to the hickory tree:

$$\mathbf{F}_a^{t+1} = \begin{cases} \mathbf{F}_a^t + d_g \times G_c \times (\mathbf{F}_h^t - \mathbf{F}_a^t), & R_1 \geq P_d \\ \text{Random location}, & R_1 < P_d \end{cases} \quad (11)$$

where t is the number of the current iterations; R_1 is a random number within $[0,1]$; P_d is the probability of predation and set as 0.1; \mathbf{F}_h^t is the squirrels in the hickory tree in the t th iteration; \mathbf{F}_a^t is the squirrels at the oak trees in the t th iteration; G_c is a sliding constant and set as 1.9; and d_g is the sliding distance [39].

2) Squirrels move from the normal trees to the oak trees:

$$\mathbf{F}_n^{t+1} = \begin{cases} \mathbf{F}_n^t + d_g \times G_c \times (\mathbf{F}_a^t - \mathbf{F}_n^t), & R_2 \geq P_d \\ \text{Random location}, & R_2 < P_d \end{cases} \quad (12)$$

where R_2 is a random probability within $[0,1]$; and \mathbf{F}_n^t is the squirrels in the normal trees at the t th iteration.

3) Squirrels move from the normal trees to the hickory tree:

$$\mathbf{F}_n^{t+1} = \begin{cases} \mathbf{F}_n^t + d_g \times G_c \times (\mathbf{F}_h^t - \mathbf{F}_n^t), & R_3 \geq P_d \\ \text{Random location}, & R_3 < P_d \end{cases} \quad (13)$$

where R_3 is a random probability within $[0,1]$.

Step 4: Seasonal changes and the termination

SSA introduces the seasonal impact in the squirrel foraging activities, to relocate squirrels in common trees terminated at wintertime. First, the seasonality constant is calculated as:

$$S_c^t = \sqrt{\sum_{k=1}^D (\mathbf{F}_{a,k}^t - \mathbf{F}_{h,k}^t)^2} \quad (14)$$

The least value of the seasonality constant (S_{\min}) is calculated as:

$$S_{\min} = \frac{10e^{-6}}{365^{t/(t_{\max}/2.5)}} \quad (15)$$

where t_{\max} is the total number of the iterations. If $S_c^t < S_{\min}$, winter comes to an end. The squirrels in the normal trees are relocated as:

$$\mathbf{F}_n^{\text{new}} = \mathbf{F}_L + \text{Levy}(n) \times (\mathbf{F}_U + \mathbf{F}_L) \quad (16)$$

where $\text{Levy}(\cdot)$ function is a mathematical model for space exploration which can be calculated as:

$$\text{Levy}(n) = 0.01 \times \frac{r_n \times \varepsilon}{|r_m|^{\frac{1}{\zeta}}} \quad (17)$$

where r_n and r_m are the two normal distributed random numbers within $[0,1]$; ζ is a constant as 1.5; and ε is calculated as:

$$\varepsilon = \left(\frac{\Gamma(1 + \zeta) \times \sin\left(\frac{\pi\zeta}{2}\right)^{\frac{1}{\zeta}}}{\Gamma\left(\frac{1+\zeta}{2}\right) \times \zeta \times 2^{\left(\frac{\zeta-1}{2}\right)}} \right)^{\frac{1}{\zeta}} \quad (18)$$

When the maximum number of the iterations is reached, the algorithm terminates. SSA can randomly generate the initial population at the initialization stage. If the initial population is not evenly distributed in the search space, the variety of the population cannot be guaranteed, which can result in insufficient optimization. Moreover, the squirrels tend to assimilate at the later stage of each iteration, which makes the SSA easily fall into local optimum.

B. Latin Hypercube Sampling Initialization

Here, this paper introduces LHS method which is a multidimensional random stratified sampling method based on Monte Carlo [44] to make the initial population of SSA evenly covered in the search space.

The search space is first divided into n independent intervals with the equal probability in LHS. Then a value from each interval with the equal probability is extracted to obtain n samples, and a vector is formed by matching the extracted n samples. Finally, the sampling set is obtained by iterative operation. The initial popu-

lation distribution produced by this method can ensure each interval to contain the same number of samples, so as to be more uniform than that of random sampling.

C. Opposition-based Learning Mechanism

In order to solve the local optimum, OBL is introduced in the SSA to select the squirrel position vector with the highest fitness value after seasonal changes so as to generate an inverse position, i.e.:

$$U_{\text{best}}^t = F_U + F_L - X_{\text{best}}^t \quad (19)$$

where U_{best}^t is the opposed squirrel position vector; and X_{best}^t is the squirrel position vector with the highest fitness after the iteration. Then, the fitness value of the opposed position is calculated and contrasted with the highest fitness value so that the squirrel with higher fitness is selected to enter the next iteration. The OBL strategy can expand the search scope, guide individual squirrels to find the global optimal solution, and avoid being trapped into local optima.

D. LOSSA Flow

Combine LHS and OBL in SSA, the proposed algorithm of LOSSA is described in Table I.

TABLE I
ALGORITHM FLOW OF LOSSA

Algorithm Pseudocode for LOSSA	
Input:	Population size n , feature dimension D , maximum number of iterations t_{max} .
Output:	Optimal feature subset GBestF.
1.	Initialize the squirrel using Latin hypercube sampling
2.	Calculate the fitness of each squirrel
3.	Sort the fitness, and declare the squirrels on hickory tree F_h , oak trees F_a and normal trees F_n
4.	while ($t < t_{\text{max}}$)
5.	for $i = 1$ to n_1 (n_1 is the total number of squirrels moving oak trees to hickory tree)
6.	Calculate F_a^{t+1} using (11)
7.	end
8.	for $i = 1$ to n_2 (n_2 is the total number of squirrels moving normal trees to oak trees)
9.	Calculate F_n^{t+1} using (12)
10.	end
11.	for $i = 1$ to n_3 (n_3 is the total number of squirrels moving normal trees to hickory tree)
12.	Calculate F_h^{t+1} using (13)
13.	end
14.	Calculate seasonal constant (S_c) using (14)
15.	Calculate the minimum value of seasonal constant (S_{min}) using (15)
16.	if $S_c < S_{\text{min}}$
17.	Relocate squirrels using (16)
18.	end
19.	Calculate the fitness of each squirrel
20.	Sort the fitness
21.	Calculate U_{best}^t using (19)
22.	if $\text{fit}(U_{\text{best}}^t) < \text{fit}(F_h^t)$
23.	$F_h^t = U_{\text{best}}^t$
24.	end
25.	end

E. AdaBoostDT Algorithm

AdaBoostDT is an adaptive boosting ensemble learning algorithm based on multiple decision tree classifiers with strong generalization capability. First, the weight D_1 of the data is initialized as:

$$D_1 = (w_{11}, w_{12}, \dots, w_{1N}) = \left(\frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N} \right) \quad (20)$$

where w_{1i} is the weight of the i th instance at the beginning; and N is the total number of the samples. After T iterations of the decision tree classifier, the decision tree classifier h_t with the lowest error rate in the current iteration is selected, and its error rate e_t is calculated as:

$$e_t = P(h_t(x_i) \neq y_i) = \sum_{i=1}^N w_{it} I(h_t(x_i) \neq y_i) \quad (21)$$

where x_i is the features of the instance; y_i represents the category label of the instance; and $I(\cdot)$ represents indicator function. According to the rate deviation, the epicycle weight a_t of the decision tree classifier is calculated as:

$$a_t = \frac{1}{2} \ln \left(\frac{1 - e_t}{e_t} \right) \quad (22)$$

Then, the weight coefficients of the training samplings are adjusted as:

$$D_{t+1} = \frac{D_t \exp(-a_t y_i h_t(x_i))}{2 \sqrt{e_t} (1 - e_t)} \quad (23)$$

Thus the decision tree classifiers are combined based on their weights to obtain the final ensemble classifier H :

$$H = \text{sign} \left(\sum_{t=1}^T a_t h_t \right) \quad (24)$$

The pseudocode of AdaBoostDT is provided in Table II.

TABLE II
ALGORITHM FLOW OF ADABOOSTDT

Algorithm Pseudocode for AdaBoostDT	
Input:	Optimal feature subset GBestF, decision tree h_t , and maximum number of the iterations T
Output:	Ensemble classifier H
1.	Initialize D_1 using (20)
2.	for t to T
3.	Train decision tree classifier h_t
4.	Calculate e_t using (21)
5.	if $e_t > 0.5$ then break
6.	Calculate a_t using (22)
7.	Calculate D_{t+1} using (23)
8.	end
9.	Get H using (24)

IV. CASE STUDIES

A. Performance Evaluation Criteria

FDIA detection is a binary anomaly detection, where the measurements can be divided into two categories:

attacked and not attacked. In this paper, the detection accuracy is represented as the ratio of correct detections. Both precision, sensitivity and F1-score are utilized to evaluate the performance and feasibility of the proposed detection method to increase the credibility of the algorithm. The confusion matrix is introduced to define the performance evaluation criteria [45].

The true positive, true negative, false positive and false negative can be seen in Table III, which are defined as follows:

- 1) The true positive represents the quantity of positives detected as positives (T_p).
- 2) The true negative represents the quantity of negatives detected as negatives (T_N).
- 3) The false positive represents the quantity of negatives detected as positives (F_p).
- 4) The false negative represents the quantity of positives detected as negatives (F_N).

TABLE III
CONFUSION MATRIX

	Detection positive	Detection negative
Real positive	Ture positive	False negative
Real negative	False positive	Ture negative

Through the confusion matrix, the four evaluation indicators can be defined as follows:

$$A = \frac{T_p + T_N}{T_p + T_N + F_p + F_N} \quad (25)$$

$$P = \frac{T_p}{T_p + F_p} \quad (26)$$

$$S = \frac{T_p}{T_p + F_N} \quad (27)$$

$$F_1 = 2 \times \frac{P \times S}{P + S} \quad (28)$$

where the precision represents the proportion of correctly detected attacked samples to all the attacked instances, and the sensitivity is the ratio of correctly detected attacked instances to the total attacked instances, reflecting the sensitivity of the classifier to the attacked samples. F1-score can balance the precision and sensitivity.

B. Experimental Setting

Most existing FDIA models are aimed to approximate direct current models with only SCADA devices. A more general FDIA model is proposed to inject false data into SCADA and PMU measurements [46]. Therefore, the FDIA model in [46] is used to attack IEEE 14-bus and IEEE 57-bus test systems. SCADA measurements include active power injection, reactive power injection, active power flow and reactive power flow, and contain measurement errors with a standard deviation of 0.02. PMU measurements include voltage phase and current phasor, and contain measurement

errors with a standard deviation of 0.001. The diagrams of the IEEE 14-bus and IEEE 57-bus test systems are presented in Figs. 3 and 4, respectively. The daily load curves for the test systems are collected from Dongguan dispatching center in Guangdong province, China in 2016, with a sampling interval of 15 minutes. Each of the two test systems produces 2480 samples, including 2000 negative samples and 480 positive samples. In order to simulate different attack scenarios, the attack model of the IEEE 14-bus test system is equipped with complete power grid topology information, while the attack model of the IEEE 57-bus test system has no such information. In the attack samples, the minimum number of attacked lines is 1, and the maximum number of attacked lines is 5. The attack strength under different number of attacked lines is calculated as:

$$\psi = \frac{1}{z-1} \sum_{i=2}^z \left| \frac{\tilde{\theta}_i - \theta_i^{\text{real}}}{\theta_i^{\text{real}}} \right| + \frac{1}{z-1} \sum_{i=2}^z \left| \frac{\tilde{V}_i - V_i^{\text{real}}}{V_i^{\text{real}}} \right| \quad (29)$$

where z is the dimension of state vector; $\tilde{\theta}_i$ is the estimated phase angel at bus i after attack; θ_i^{real} is the real phase angel; \tilde{V}_i is the estimated voltage magnitude at bus i ; and V_i^{real} is the real voltage magnitude. Table IV lists the results of ψ with different number of attacked lines in the IEEE 14-bus and IEEE 57-bus test systems.

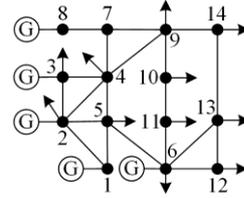


Fig. 3. The diagram of the IEEE 14-bus test system.

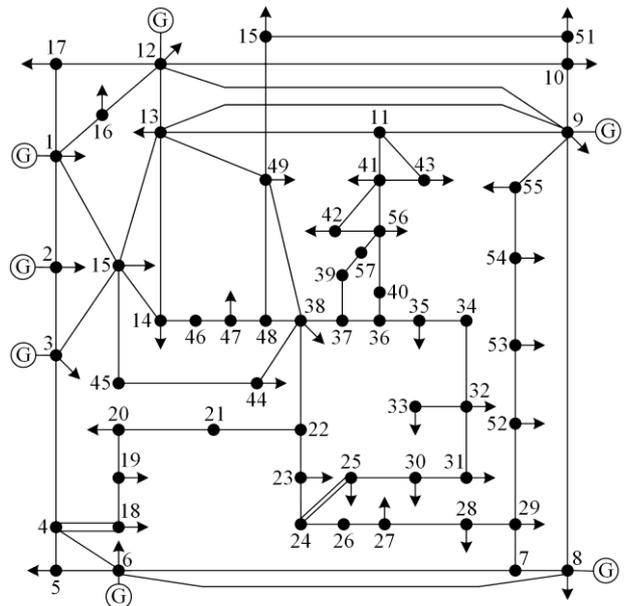


Fig. 4. The diagram of the IEEE 57-bus test system.

TABLE IV
THE VALUES OF ψ UNDER ATTACKED LINES

Test system	Number of attacked lines				
	1	2	3	4	5
14-bus	0.0536	0.1628	0.1932	0.2438	0.3180
57-bus	0.4202	0.5716	0.7140	1.4735	1.4835

Table V lists the average amount of features in the optimal feature subset of each label in the two datasets of the IEEE 14-bus and IEEE 57-bus test systems.

TABLE V
NUMBER OF FEATURES IN THE OPTIMAL FEATURE SUBSET

Test system	Features	Selected Features
14-bus	104	53
57-bus	419	209

For the model training, the ratio of training set and testing set is 7:3. The number of the weak classifiers for AdaBoostDT is set as 100. In the IEEE 14-bus and IEEE 57-bus test systems under FDIA, AdaBoostDT, random forest (RF) and extreme learning machine (ELM) algorithms are applied for comparing the accuracy, precision, sensitivity and F1-score with the proposed method.

C. Simulation Experiment of the IEEE 14-bus Test System

First, the traditional FDIA model with the complete topology information of the CPPS is considered. Here, the FDIA in the IEEE 14-bus test system is an attack with complete known information of the power system. The FDIA detection accuracies of LOSSA-AdaBoostDT, AdaBoostDT, RF and ELM are demonstrated in Fig. 5. The state variables represent the phase angle and voltage magnitude of each bus except the balancing node, which are arranged according to the order of the bus in Fig. 5.

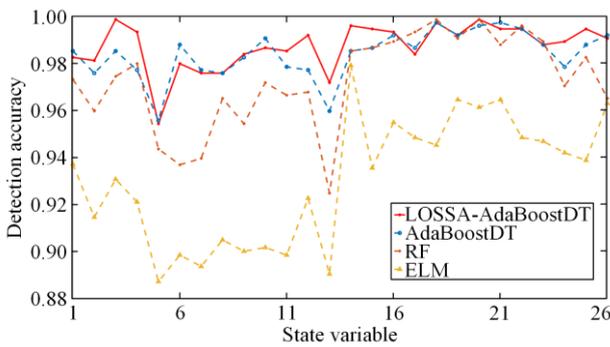


Fig. 5. The detection accuracies of the IEEE 14-bus test system.

As depicted in Fig. 5, the RF and ELM algorithms have low detection accuracies and large fluctuations, with average accuracies of 97.27% and 93.04%, while the minimum accuracies are only 92.47% and 88.71%, respectively. In comparison, the accuracy of AdaBoostDT is significantly higher than those of RF and ELM, and it is more stable for different node states. Compared with the AdaBoost algorithm, the detection

accuracy of LOSSA-AdaBoostDT with additional feature selection can be further improved, and the average detection accuracy reaches 98.72%. It demonstrates that the average accuracy of LOSSA-AdaBoostDT is better than those of other methods in the same attack model.

Table VI lists the precision, sensitivity and F1-score of the four algorithms in the IEEE 14-bus test system. The average detection precision of LOSSA-AdaBoostDT is 1.68%, 3.05% and 3.06% higher than AdaBoostDT, RF and ELM, respectively, and can reach up to 98.24%. The sensitivity and F1-score of the LOSSA-AdaBoostDT are 98.51% and 98.38%, respectively, being higher than other tree algorithms. Hence, the sensitivity of the proposed algorithm to the attacked samples and the ability to correctly detect the attacked samples are superior to the other three algorithms.

TABLE VI
DETECTION EVALUATION IN THE IEEE 14-BUS TEST SYSTEM

Test system	Method	Index		
		Precision	Sensitivity	F1-score
14-bus	LOSSA-AdaBoostDT	0.9824	0.9851	0.9838
	AdaBoostDT	0.9656	0.9817	0.9737
	RF	0.9519	0.9779	0.9649
	ELM	0.9518	0.8596	0.9057

D. Simulation Experiment of the IEEE 57-bus Test System

FDIA detection is further applied on the IEEE 57-bus test system without complete power system information. The detection accuracies of the test system with diverse algorithms are demonstrated in Fig. 6.

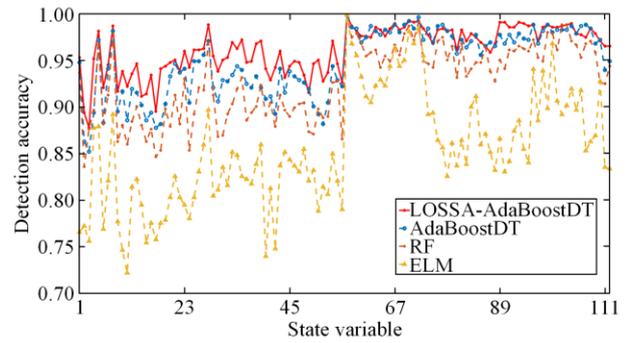


Fig. 6. The detection accuracies of the IEEE 57-bus test system.

In Fig. 6, the average accuracies of the LOSSA-AdaBoostDT, AdaBoostDT, RF and ELM detection methods are 96.21%, 94.82%, 92.77% and 85.53% respectively. Compared with the IEEE 14-bus system, the detection accuracies of AdaBoostDT, RF and ELM decrease significantly in the IEEE 57-bus test system. The data dimension of the IEEE 57-bus test system is 419, which is much higher than the 104 dimensions in the IEEE 14-bus test system. Specifically, the high dimensionality of the dataset leads to the performance degradation of the classifier. After the introduction of LOSSA

feature selection, the dimension of the IEEE 57-bus test system is reduced to 209, which can maintain good detection performance of the classifier. Table VII lists the precision, sensitivity and F1-score of the different algorithms in the IEEE 57-bus test system. It is shown that the average detection precision of LOSSA-AdaBoostDT is up to 92.76%, and is 1.78%, 8.05% and 8.56% higher than AdaBoostDT, RF and ELM, respectively. In addition, LOSSA-AdaBoostDT results in the sensitivity and F1-score index of 95.18% and 93.97%, also being higher than the other tree algorithms.

TABLE VII
DETECTION EVALUATION IN THE IEEE 57-BUS TEST SYSTEM

Test system	Method	Index		
		Precision	Sensitivity	F1-score
57-bus	LOSSA-AdaBoostDT	0.9276	0.9518	0.9397
	AdaBoostDT	0.9098	0.9229	0.9164
	RF	0.8471	0.9153	0.8812
	ELM	0.8420	0.6177	0.7299

E. Receiver Operating Characteristic Analysis

Receiver operating characteristics [47] (ROC) is used to describe the relationship between true positive rate (TPR) and false positive rate (FPR). The ROC curves of the IEEE 14-bus and 57-bus test systems are depicted in Fig. 7. The area covered by ROC curve and the coordinate axis is the area under curve (AUC) [47]. Since the function of a classifier is to classify samples into positive and negative classes, AUC can represent the classifier's ability to distinguish positive and negative samples. A larger AUC indicates higher detection performance of the algorithm, and the ideal algorithm has an AUC value of 1. In Fig. 7, the TPR and FPR of the IEEE 14-bus test system are 0.9861 and 0.0066, respectively, while the TPR and FPR of the IEEE 57-bus test system are 0.9579 and 0.0325, respectively. The AUCs of both examples are close to 1, further demonstrating that the proposed method has high FDIA detection accuracy.

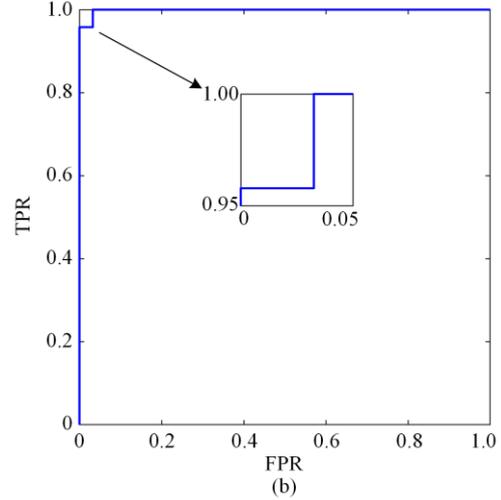
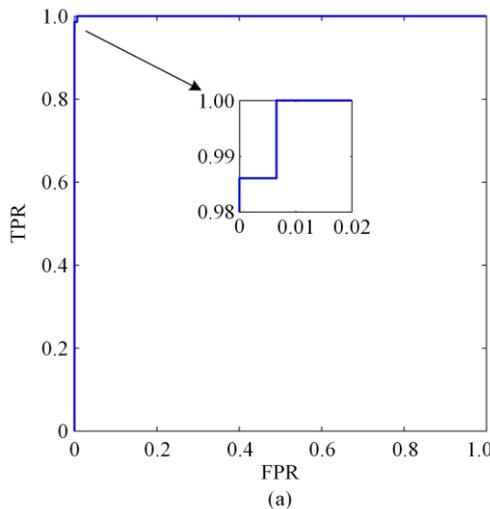


Fig. 7. ROC curves of the IEEE test systems. (a) IEEE 14-bus test system. (b) IEEE 57-bus test system.

V. CONCLUSION

In this paper, a LOSSA-AdaBoostDT method is proposed to address the weak generalization ability of the single classifier in machine learning-based FDIA detection methods for CPPS, and the problem of the original physical meaning of the measurement features altering during dimension reduction. The proposed method adopts the AdaBoostDT classifier with strong generalization ability to detect FDIA more quickly and accurately. At the same time, LOSSA is applied with the combining of LHS and OBL to avoid the local optimum of SSA and obtain the optimal feature subset to improve the detection performance of the classifier.

Simulation experiments in the IEEE 14-bus and 57-bus test systems are performed to prove the effectiveness of the LOSSA-AdaBoostDT. Compared with AdaBoostDT, RF and ELM algorithms, LOSSA-AdaBoostDT has better detection performance, and the detection accuracy of IEEE 14-bus and IEEE 57-bus test systems reach 98.72% and 96.21%, respectively. Further study will investigate new methods to eliminate and recover the contaminated states in the CPPS so as to improve the stability of the CPPS.

ACKNOWLEDGMENT

The authors thank Hubei Provincial Key Laboratory for Operation and Control of Cascaded Hydropower Station, China Three Gorges University for assistance in theoretical analysis and experiments.

AUTHORS' CONTRIBUTIONS

Lei Xi: methodology, formal analysis, software, visualization, and writing original draft. Xilong Tian: writing, reviewing, editing, and validation. Miao He: data collection and management. Chen Cheng: conceptualization and supervision. All authors read and approved the final manuscript.

FUNDING

This work is supported by the National Natural Science Foundation of China (No. 52477104).

AVAILABILITY OF DATA AND MATERIALS

Not applicable.

DECLARATIONS

Competing interests: The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this article.

AUTHORS' INFORMATION

Lei Xi received the M.S. degree in control theory and control engineering from the Harbin University of Science and Technology, Harbin, China, in 2009, and Ph.D. degree in electrical engineering from the South China University of Technology, Guangzhou, China, in 2013. He is currently a full professor with the College of Electrical Engineering and New Energy, China Three Gorges University, Yichang, China. His research interests include load frequency control, artificial intelligence techniques, automatic generation control, and network attack & defense.

Xilong Tian received the bachelor's degree in information management and information system from Hebei University of Technology, Tianjin, China, in 2019, and the M.S. degree in electrical engineering from China Three Gorges University, Yichang, China, in 2024. Now he is with Qianjiang Power Supply Company, State Grid Hubei Electric Power Co., Ltd., Qianjiang, China. His research interests include cyber-physical power system and network attack and defense.

Miao He received the bachelor's degree in electrical engineering from Xi'an University of Technology, Xi'an, China, in 2020, and the M.S. degree in electrical engineering from China Three Gorges University, Yichang, China, in 2023. Her research interests include cyber-physical power system and network attack & defense.

Chen Cheng received the bachelor's degree in electrical engineering from Qingdao University of Technology, Qingdao, China, in 2021, and the M.S. degree in electrical engineering from China Three Gorges University, Yichang, China, in 2024. Her research interests include cyber-physical power system and network attack & defense.

REFERENCES

- [1] B. Yan, Z. Jiang, and P. Yao, "Game theory based optimal defensive resources allocation with incomplete information in cyber-physical power systems against false data injection attacks," *Protection and Control of Modern Power Systems*, vol. 9, no. 2, pp. 115-127, Mar. 2024.
- [2] W. Hao, Q. Yang, and Z. Li *et al.*, "Multi-scale traffic aware cybersecurity situational awareness online model for intelligent power substation communication network," *IEEE Internet of Things Journal*, vol. 10, no. 2, pp. 1666-1681, Jan. 2023.
- [3] Y. Wu, Y. Ru, and Z. Lin *et al.*, "Research on cyber-attacks and defensive measures of power communication network," *IEEE Internet of Things Journal*, vol. 10, no. 9, pp. 7613-7635, May 2023.
- [4] Y. Wu, H. Xu, and M. Ni, "Defensive resource allocation method for improving survivability of communication and information system in CPPS against cyber-attacks," *Journal of Modern Power Systems and Clean Energy*, vol. 8, no. 4, pp. 750-759, Jul. 2020.
- [5] Y. Xu, "A review of cyber security risks of power systems: from static to dynamic false data attacks," *Protection and Control of Modern Power Systems*, vol. 5, no. 3, pp. 1-12, Jul. 2020.
- [6] B. Chen, Z. Yang, and Y. Zhang *et al.*, "Risk assessment of cyber-attacks on power grids considering the characteristics of attack behaviors," *IEEE Access*, vol. 8, pp. 148221-148344, Aug. 2020.
- [7] Y. Liu, H. B. Gooi, and Y. Li *et al.*, "A secure distributed transactive energy management scheme for multiple interconnected microgrids considering misbehaviors," *IEEE Transactions on Smart Grid*, vol. 10, no. 6, pp. 5975-5986, Nov. 2019.
- [8] W. Hao, P. Yao, and T. Yang *et al.*, "Industrial cyber-physical system defense resource allocation using distributed anomaly detection," *IEEE Internet of Things Journal*, vol. 9, no. 22, pp. 22304-22314, Nov. 2022.
- [9] Y. Liu, Y. Li, and Y. Wang *et al.*, "Robust and resilient distributed optimal frequency control for microgrids against cyber attacks," *IEEE Transactions on Industrial Informatics*, vol. 18, no. 1, pp. 375-386, Jan. 2022.
- [10] Y. Liu, H. Xin, and Z. Qu *et al.*, "An attack-resilient cooperative control strategy of multiple distributed generators in distribution networks," *IEEE Transactions on Smart Grid*, vol. 7, no. 6, pp. 2923-2932, Nov. 2016.
- [11] Y. Wu, J. Chen, and Y. Ru *et al.*, "Research on power communication network planning based on information transmission reachability against cyber-attacks," *IEEE Systems Journal*, vol. 15, no. 2, pp. 2883-2894, Jun. 2021.
- [12] Y. Chakhchoukh and H. Ishii, "Enhancing robustness to cyber-attacks in power systems through multiple least trimmed squares state estimations," *IEEE Transactions on Power Systems*, vol. 31, no. 6, pp. 4395-4405, Nov. 2016.
- [13] B. Chen, H. Li, and B. Zhou, "Real-time identification of false data injection attacks: A novel dynamic-static parallel state estimation based mechanism," *IEEE Access*, vol. 7, pp. 95812-95824, Jul. 2019.
- [14] M. Jorjani, H. Seifi, and A. Y. Varjani, "A graph theory-based approach to detect false data injection attacks in power system AC state estimation," *IEEE Transactions on Industrial Informatics*, vol. 17, no. 4, pp. 2465-2475, Apr. 2021.
- [15] H. M. Khalid and J. C. Peng, "Immunity toward data-injection attacks using multi-sensor track fusion-based model prediction," *IEEE Transactions on Smart Grid*, vol. 8, no. 2, pp. 697-707, Mar. 2017.
- [16] S. Li, Y. Yilmaz, and X. Wang, "Quickest detection of false data injection attack in wide-area smart grids," *IEEE Transactions on Smart Grid*, vol. 6, no. 6, pp. 2725-2735, Nov. 2015.
- [17] J. Zhao, G. Zhang, and M. L. Scala *et al.*, "Short-term state forecasting-aided method for detection of smart grid

- general false data injection attacks,” *IEEE Transactions on Smart Grid*, vol. 8, no. 4, pp. 1580-1590, Jul. 2017.
- [18] S. Basumallik, R. Ma, and S. Eftekharnjad, “Packet-data anomaly detection in PMU-based state estimator using convolutional neural network,” *International Journal of Electrical Power & Energy Systems*, vol. 107, pp. 690-702, May 2019.
- [19] L. Xi, L. Dong, and C. Cheng *et al.*, “Location detection of a false data injection attack in acyber-physical power system based on a hybrid chimp optimized extreme learning machine,” *Power System Protection and Control*, vol. 52, no. 14, pp. 46-58, Jul. 2024. (in Chinese)
- [20] B. Chen, Q. Wu, and M. Li *et al.*, “Detection of false data injection attacks on power systems using graph edge-conditioned convolutional networks,” *Protection and Control of Modern Power Systems*, vol. 8, no. 2, pp. 1-12, Apr. 2023.
- [21] A. Sayghe, Y. Hu, and I. Zografopoulos *et al.*, “Survey of machine learning methods for detecting false data injection attacks in power systems,” *IET Smart Grid*, vol. 3, no. 5, pp. 581-595, Oct. 2020.
- [22] X. Huang, Z. Qin, and M. Xie *et al.*, “Defense of massive false data injection attack via sparse attack points considering uncertain topological changes,” *Journal of Modern Power Systems and Clean Energy*, vol. 6, pp. 1588-1598, Sept. 2022.
- [23] A. Aburomman and M. Reaz, “A survey of intrusion detection systems based on ensemble and hybrid classifiers,” *Computers & Security*, vol. 65, pp. 135-152, Mar. 2017.
- [24] X. Dong, Z. Yu, and W. Cao, *et al.*, “A survey on ensemble learning,” *Frontiers of Computer Science*, vol. 14, no. 2, pp. 241-258, Apr. 2020.
- [25] P. K. Jena, S. Ghosh, and E. Koley, *et al.*, “An ensemble classifier based scheme for detection of false data attacks aiming at disruption of electricity market operation,” *Journal of Network and Systems Management*, vol. 29, no. 4, pp. 43, Jun. 2021.
- [26] X. Lu, J. Jing, and Y. Wu, “False data injection attack location detection based on classification method in smart grid,” in *2020 2nd International Conference on Artificial Intelligence and Advanced Manufacture (AIAM)*, Manchester, UK, Oct. 2020, pp. 133-136.
- [27] W. Xue and T. Wu, “Active learning-based XGBoost for cyber physical system against generic AC false data injection attacks,” *IEEE Access*, vol. 8, pp. 144575-144584, Aug. 2020.
- [28] P. Zhang and Z. Yang, “A novel AdaBoost framework with robust threshold and structural optimization,” *IEEE Transactions on Cybernetics*, vol. 48, no. 1, pp. 64-76, Jan. 2018.
- [29] S. Wan, X. Li, and Y. Yin *et al.*, “Milling chatter detection by multi-feature fusion and Adaboost-SVM,” *Mechanical Systems and Signal Processing*, vol. 156, Jul. 2021.
- [30] J. Si, S. Tu, and X. Fan, “Fault diagnosis of high-voltage circuit breaker based on SO-PAA-GAF and AdaBoost ensemble learning,” *Power System Protection and Control*, vol. 52, no. 3, pp.152-160, Feb. 2024. (in Chinese)
- [31] Y. Liu, C. Zhao, and H. Liang *et al.*, “A rotor fault diagnosis method based on BP-Adaboost weighted by non-fuzzy solution coefficients,” *Measurement*, vol. 196, Jun. 2022.
- [32] W. You, K. Shen, and N. Yang, “Research on electricity theft detection based on AdaBoost ensemble learning,” *Power System Protection and Control*, vol. 48, no. 19, pp.151-159, Oct. 2020. (in Chinese)
- [33] A. Jiang, H. Wei, and J. Deng *et al.*, “Cloud-edge cooperative model and closed-loop control strategy for the price response of large-scale air conditioners considering data packet dropouts,” *IEEE Transactions on Smart Grid*, vol. 11, no. 5, pp. 4201-4211, Sept. 2020.
- [34] G. Chao, Y. Luo, and W. Ding, “Recent advances in supervised dimension reduction: a survey,” *Machine Learning and Knowledge Extraction*, vol. 1, no. 1, pp. 341-358, Jan. 2019.
- [35] Y. Zhang, J. Wang, and B. Chen, “Detecting false data injection attacks in smart grids: A semi-supervised deep learning approach,” *IEEE Transactions on Smart Grid*, vol. 12, no. 1, pp. 623-634, Jan. 2021.
- [36] A. Parizad and C. Hatziaioniu, “Cyber-attack detection using principal component analysis and noisy clustering algorithms: a collaborative machine learning-based framework,” *IEEE Transactions on Smart Grid*, vol. 13, pp. 4848-4861, Nov. 2022.
- [37] S. Ahmed, Y. Lee, and S. H. Hyun *et al.*, “Feature selection-based detection of covert cyber deception assaults in smart grid communications networks using machine learning,” *IEEE Access*, vol. 6, pp. 27518-27529, May 2018.
- [38] B. H. Nguyen, B. Xue, and M. Zhang, “A survey on swarm intelligence approaches to feature selection in data mining,” *Swarm and Evolutionary Computation*, vol. 54, May 2020.
- [39] M. Jain, V. Singh, and A. Rani, “A novel nature-inspired algorithm for optimization: Squirrel search algorithm,” *Swarm and Evolutionary Computation*, vol. 44, pp. 148-175, Feb. 2019.
- [40] Y. Wang and T. Du, “An improved squirrel search algorithm for global function optimization,” *Algorithms*, vol. 12, no. 4, pp. 80, Apr. 2019.
- [41] K. Sharma and M. K. Trivedi, “Latin hypercube sampling-based NSGA-III optimization model for multi-mode resource constrained time-cost-quality-safety trade-off in construction projects,” *International Journal of Construction Management*, pp. 1-11, Nov. 2020.
- [42] M. Tubishat, N. Idris, and L. Shuib *et al.*, “Improved salp swarm algorithm based on opposition based learning and novel local search algorithm for feature selection,” *Expert Systems with Applications*, vol. 145, May 2020.
- [43] M. Ashrafuzzaman, S. Das, and Y. Chakhchoukh *et al.*, “Detecting stealthy false data injection attacks in the smart grid using ensemble-based machine learning,” *Computers & Security*, vol. 97, Oct. 2020.
- [44] G. Bayar and G. Hambarci, “Improving measurement accuracy of indoor positioning system of a Mecanum wheeled mobile robot using Monte Carlo-Latin hypercube sampling based machine learning algorithm,” *Journal of the Franklin Institute*, vol. 360, no. 17, pp. 13994-14021, Aug. 2022.
- [45] W. Hao, T. Yang, and Q. Yang, “Hybrid statistical-machine learning for real-time anomaly detection in industrial cyber-physical systems,” *IEEE Transactions on Automation Science and Engineering*, vol. 20, no. 1, pp. 32-46, Jan. 2022.
- [46] T. Wu, W. Xue, and W. Wang *et al.*, “Extreme learning machine-based state reconstruction for automatic attack filtering in cyber physical power system,” *IEEE Transactions on Industrial Informatics*, vol. 17, no. 3, pp. 1892-1904, Mar. 2021.
- [47] W. Lei, L. G. A. Alves, and L. A. N. Amaral, “Forecasting the evolution of fast-changing transportation networks using machine learning,” *Nature Communications*, vol. 13, no. 1, pp. 4252, Jul. 2022.