

ConAML: Constrained Adversarial Machine Learning for Cyber-Physical Systems

Jiangnan Li
jli103@vols.utk.edu
University of Tennessee, Knoxville

Yingyuan Yang
yyang260@uis.edu
University of Illinois Springfield

Jinyuan Stella Sun
jysun@utk.edu
University of Tennessee, Knoxville

Kevin Tomsovic
tomsovic@utk.edu
University of Tennessee, Knoxville

Hairong Qi
hqi@utk.edu
University of Tennessee, Knoxville

ABSTRACT

Recent research demonstrated that the superficially well-trained machine learning (ML) models are highly vulnerable to adversarial examples. As ML techniques are becoming a popular solution for cyber-physical systems (CPSs) applications in research literatures, the security of these applications is of concern. However, current studies on adversarial machine learning (AML) mainly focus on pure cyberspace domains. The risks the adversarial examples can bring to the CPS applications have not been well investigated. In particular, due to the distributed property of data sources and the inherent physical constraints imposed by CPSs, the widely-used threat models and the state-of-the-art AML algorithms in previous cyberspace research become infeasible.

We study the potential vulnerabilities of ML applied in CPSs by proposing Constrained Adversarial Machine Learning (ConAML), which generates adversarial examples that satisfy the intrinsic constraints of the physical systems. We first summarize the difference between AML in CPSs and AML in existing cyberspace systems and propose a general threat model for ConAML. We then design a best-effort search algorithm to iteratively generate adversarial examples with linear physical constraints. We evaluate our algorithms with simulations of two typical CPSs, the power grids and the water treatment system. The results show that our ConAML algorithms can effectively generate adversarial examples which significantly decrease the performance of the ML models even under practical constraints.

KEYWORDS

adversarial machine learning; cyber-physical system; intrusion detection

1 INTRODUCTION

Machine learning (ML) has shown promising performance in many real-world applications, such as image classification [20], speech recognition [18], and malware detection [49]. In recent years, motivated by the promotion of cutting-edge communication and computational technologies, there is a trend to adopt ML in various cyber-physical system (CPS) applications, such as data center thermal management [29], agriculture ecosystem management [9], power grid attack detection [37], and industrial control system anomaly detection [24].

Recent research has demonstrated that the superficially well-trained ML models are highly vulnerable to adversarial examples

[10, 17, 26, 34, 35, 40, 42]. In particular, adversarial machine learning (AML) technologies enable attackers to deceive ML models with well-crafted adversarial examples by adding small perturbations to legitimate inputs. As CPSs have become synonymous to security-critical infrastructures such as the power grid, nuclear systems, avionics, and transportation systems, such vulnerabilities can be exploited leading to devastating consequences.

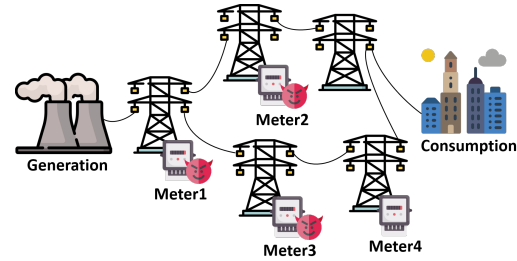


Figure 1: A CPS example (power grids).

AML research has received considerable attention in artificial intelligence (AI) communities and it mainly focuses on computational applications such as computer vision. However, it is not applicable to CPSs because the inherent properties of CPSs render the widely-used threat models and AML algorithms in previous research infeasible. The existing AML research makes common assumptions on the attacker's knowledge and the adversarial examples. In most AML research, the attacker is assumed to have full knowledge of the ML inputs and these features are assumed to be mutually independent. For example, in computer vision [17], the attacker is assumed to know all the values of pixels of an image and there is no strict dependency among the pixels. However, this is not realistic for attacks targeting CPSs. CPSs are usually large and complex systems whose data sources are heterogeneous and geographically distributed. The attacker may compromise a subset of sensors and modify their measurement data. Generally, for the uncompromised data sources, the attacker cannot even know the measurements, let alone making modifications. Furthermore, for robustness and resilience reasons, CPSs usually employ redundant data sources and incorporate faulty data detection mechanisms. For example, in the power grid, redundant phasor measurement units (PMUs) are deployed in the field to measure frequency and phase angle, and residue-based bad data detection is employed to detect and recover from faulty data for state estimation [45]. Therefore, the

features of ML applications in CPS are not only dependent but also subject to the physical constraints of the system. A simple example of constraints is shown in Figure 1. All three meters are measuring the electric current (Ampere) data. If an attacker compromises *Meter1*, *Meter2*, and *Meter3*, no matter what modification the attacker makes to the measurements, the compromised measurement of *Meter1* should always be the sum of that of *Meter2* and *Meter3* due to Kirchhoff’s laws. Otherwise, the crafted measurements will be detected by the bad data detection mechanism and obviously anomalous to the power system operators. In addition to distributed data sources and physical constraints, sensors in real-world CPSs are generally configured to collect data with a specific sampling rate. A valid adversarial attack needs to be finished within the CPS’ sampling period.

The intrinsic properties of CPS pose stringent requirements for the adversarial attackers. The attacker is now required to overcome:

- **Knowledge constraint:** No access to the ML models and the measurement values of uncompromised sensors.
- **Physical constraint:** The adversarial examples need to meet the physical constraints defined by the system.
- **Time constraint:** Attacks need to be completed within a sample period of the sensors.

to launch an effective attack that deceives the ML applications deployed in CPSs. However, in this paper, we show that the ML applications in CPSs are susceptible to handcrafted adversarial examples even though such systems naturally pose a greater barrier for the attacker.

In this paper, we propose constrained adversarial machine learning (**ConAML**), a general AML framework that incorporates the above constraints of CPSs. We firstly design a universal adversarial measurement algorithm to solve the knowledge constraint. After that, without loss of generality, we present a practical best-effort search algorithm to effectively generate adversarial examples under linear physical constraints which are one of the most common constraints in real-world CPS applications, such as power grids [33] and water pipelines [16]. Meanwhile, we set the maximum iteration number to control the time cost of the attack. We implement our algorithms with ML models used in two CPSs and mainly focus on neural networks due to its transferability. Our main contributions are summarized as follows:

- We highlight the potential vulnerability of deploying ML in CPSs, analyze the different requirements for AML applied in CPSs with regard to the general computational applications, and present a practical threat model for AML in CPSs.
- We formulate the mathematical model of ConAML by incorporating the physical constraints of the underlying system.
- We proposed ConAML, an AML framework that contains a series of AML algorithms to generate adversarial examples under the corresponding constraints.
- We assess our algorithms with two typical CPSs, the power grids and water treatment system, where ML are intensively investigated for attack detection in the research literature [1, 3, 7, 12, 13, 21–24, 36, 37, 44, 46]. The evaluation results show that the adversarial examples generated by our algorithms can effectively bypass the ML-powered attack detection systems in the two CPSs.

Related research is discussed in Section 2. We analyze the properties of AML in CPSs and give the mathematical definition and the threat model in Section 3. Section 4 presents the algorithm design. Section 5 uses two CPSs as proofs of concept to carry our experiments. Discussions and future work are given in Section 6. Section 7 concludes the paper.

2 RELATED WORK

AML of deep neural network (DNN) was discovered by Szegedy *et al.* [42] in 2013. They found that a DNN used for image classification can be fooled by adding a hardly perceptible perturbation to the legitimate image. The same perturbation can cause a different DNN to misclassify the same image even when the DNN has a different structure and is trained with a different dataset, which is referred to as the transferability property of adversarial examples. In 2015, Goodfellow *et al.* [17] proposed the Fast Gradient Sign Method (FGSM), an efficient algorithm to generate adversarial examples. The Fast Gradient Value (FGV) method by Rozsa *et al.* [40] is a variant of FGSM and utilizes the raw gradient instead of the sign values. Moosavi-Dezfooli *et al.* presented *DeepFool* to iteratively search for the closest distance between the original input and the decision boundary [34]. Single-step attacks have better transferability but can be easily defended [26]. Therefore, multi-steps methods, such as iterative methods [26] and momentum-based methods [10], are presented. The above methods generate individual adversarial examples for each input. In 2017, Moosavi-Dezfooli *et al.* designed universal adversarial perturbations to generate perturbations regardless of the ML model inputs [35].

Research on AML applications continues growing rapidly. Sharif *et al.* launched adversarial attacks to a face-recognition system and achieved a notable result [41]. Grossee *et al.* constructed adversarial attacks against Android malware detection models [19]. In 2014, Laskov *et al.* developed a taxonomy for practical adversarial attacks based on the attackers’ capability and launched evasion attacks to *PDFRATE*, a real-world online machine learning system to detect malicious PDF malware [39]. In 2018, Li *et al.* presented *TEXTBUGGER*, a framework to generate adversarial text against deep learning-based text understanding (DLTU) systems and achieved state-of-the-art attack performance [28].

AML techniques that involve the physical domain are drawing more and more attention. Kurakin *et al.* presented that ML models are vulnerable to adversarial examples in physical world scenarios by feeding a phone camera captured adversarial image to an ImageNet classifier [25]. In 2016, Carlini *et al.* presented that well-crafted voice commands which are unintelligible to human listeners, can be interpreted as commands by voice controllable systems [5]. [43] and [32] investigated the security of ML models used in autonomous driving cars. In 2018, [15] showed that an attacker can generate adversarial examples by modifying a portion of measurements in CPSs, and presented an anomaly detection model where each sensor’s reading is predicted as a function of other sensors’ readings. After that, Erba *et al.* also studied the AML in CPS and consider the physical constraints [11]. They employed an autoencoder that is trained on normal system data to reconstruct the bad inputs to match the physical behavior. However, both [15] and [11] allow

the attacker to know all the measurements which may be impractical in real-world attacks. Meanwhile, the generated adversarial examples of [11] may still violate the physical constraints.

More related work on adversarial attacks, including the adversarial example generation and applications, can be found in [48].

3 SYSTEM AND THREAT MODEL

3.1 ML-Assisted CPSs

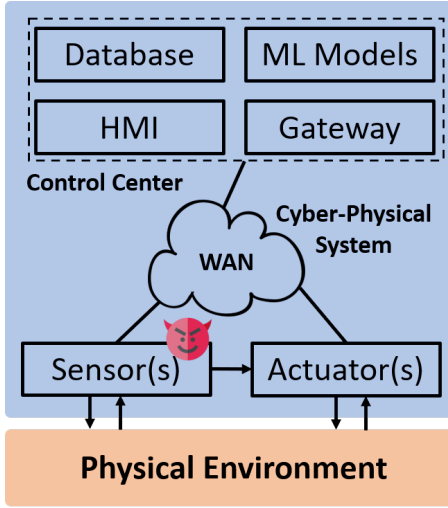


Figure 2: Machine learning-assisted CPS architecture.

Generally, a CPS can be simplified as a system that consists of four parts, namely sensors, actuators, the communication network, and the control center [7], as shown in Figure 2. The sensors measure and quantify the data from the physical environment, and send the measurement data to the control center through the communication network. In practice, the raw measurement data will be filtered and processed by the gateway according to the error checking mechanism whose rules are defined by human experts based on the properties of the physical system. Measurement data that violates the physically defined rules will be removed.

Similar to [11], we consider the scenario that the control center utilizes ML model(s) to make decisions (classification) based on the filtered measurement data from the gateway directly, and the features used to train the ML models are the measurements of sensors respectively. The target of the attacker will be deceiving the ML model(s) in CPSs to output wrong (classification) results without being detected by the gateway by adding perturbations to the measurements of the compromised sensors.

3.2 Threat Model

Adversarial attacks can be classified according to the attacker’s capability and attack goals [6, 39, 48]. In this work, we consider the integrity attack that the attacker generates adversarial perturbations to the ML inputs to deceive the ML model to make incorrect classification outputs.

There are several inherent properties of CPS that pose specific requirements for adversarial attacks. First, in CPS, ML models are

usually placed in the control centers and other centralized locations which employ comprehensive and advanced security measures such as air-gapped networks. It is highly unlikely for the attacker to have access to the models and a black-box attack should be considered. Second, we assume that the attacker cannot access the training dataset for the same reason as above, but has access to an alternative dataset such as historical data that follows a similar distribution to train their models. It is possible for the attacker to obtain historical data in practice, for instance, temperature data for load forecasting, earthquake sensor data, flood water flow data, and traffic flow data, since these data are usually published or shared among multiple parties.

To launch adversarial attacks, the attacker is assumed to compromise a certain number of sensors, and can freely eavesdrop and modify their measurement data. These sensors are deployed in the wild and their security is hard to guarantee. In real attack scenarios, this can be implemented by either directly compromising the sensors, such as device intrusion or attacking the communication network, such as man-in-the-middle attacks. However, due to the vastly distributed nature of sensors in CPS, it is only reasonable for the attacker to compromise a subset of the data sources but not all of them. For the uncompromised sensors, the attacker can neither know their measurement values nor make modifications. This constraint indicates that the attacker has limited knowledge of the ML inputs.

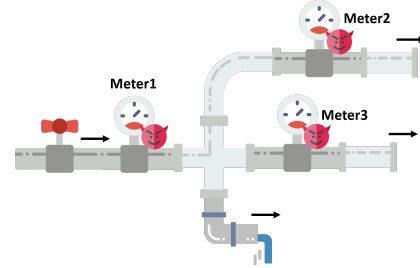


Figure 3: A CPS example (water pipelines).

Meanwhile, the attacker is further required to generate adversarial examples that meet the constraints imposed by the physical laws and system topology and evade any built-in detection mechanisms in the system. Specifically, since they are very common in real-world CPSs, we will mainly focus on linear constraints in this paper, including both linear equality constraints and linear inequality constraints. An example of the linear inequality constraint is shown in Figure 3. All the meters in Figure 3 are measuring water flow which follows the arrows’ direction. If an attacker wants to defraud the anomaly detection ML model of a water treatment system by modifying the meters’ readings, the adversarial measurement of *Meter1* should always be larger than the sum of *Meter2* and *Meter3* due to the physical structure of the pipelines. Otherwise, the poisoned inputs will be obviously anomalous to the victim (system operator) and detected automatically by the error checking mechanisms. In practice, many of the linear constraints can be explicitly abstracted by the attacker if she/he obtained enough measurement data by observing the compromised sensors. Meanwhile, the practical CPSs

usually have built-in tolerance for noise and normal fluctuation in the measurements so that the approximately estimated constraints will still be effective for the adversarial attackers. Therefore, we assume that the attacker know the linear constraints among the compromised measurements. We discuss the nonlinear equality constraints at **Appendix A**.

The real-world CPSs, such as the Supervisory Control and Data Acquisition (SCADA), will have a constant measurement sampling rate (frequency) configured for their sensors. The attacker who targets CPSs' ML applications is then required to generate a valid adversarial example within a measurement sampling period.

We summarize the threat model as follows:

- We assume the attacker has no access to the system operator's trained model in the control center, including the hyper-parameters and the related dataset. However, the attacker has an alternative dataset as an approximation of the defender's (system operator's) training dataset to train his/her ML models.
- The attacker can compromise a subset of sensors in the CPS and make modifications to their measurement data. However, the attack can neither know nor modify the measurements of uncompromised sensors.
- The attacker can know the linear constraints of the measurements imposed by the physical system.

3.3 Physical Constraint Mathematical Representation

In this subsection, we present the mathematical definition of the physical linear constraints of the ML inputs and represent the AML as a constrained optimization problem.

3.3.1 Notations. To simplify the mathematical representation, we will use $A_B = [a_{b_0}, a_{b_1}, \dots, a_{b_{n-1}}]$ to denote a sampled vector of $A = [a_0, a_1, \dots, a_{m-1}]$ according to B , where $B = [b_0, b_1, \dots, b_{n-1}]$ is a vector of sampling index. For example, if $A = [a, b, c, d, e]$ and $B = [0, 2, 4]$, we have $A_B = [a, c, e]$.

We assume there are totally d sensors in a CPS, and each sensor's measurement is a feature of the ML model f_θ in the control center. We use $S = [s_0, s_1, \dots, s_{d-1}]^T$ and $M = [m_0, m_1, \dots, m_{d-1}]^T$ to denote all the sensors and their measurements respectively. The attacker compromised r sensors in the CPS and $C = [c_0, c_1, \dots, c_{r-1}]$ denotes the index vector of the compromised sensors. Obviously, we have $\|C\| = r$ and $0 < r \leq d$. Meanwhile, the uncompromised sensors' indexes are denoted as $U = [u_0, u_1, \dots, u_{d-r-1}]$ ($\|U\| = d - r$).

$\Delta = [\delta_0, \delta_1, \dots, \delta_{d-1}]^T$ is the adversarial perturbation to be added to M . However, the attacker can only inject $\Delta_C = [\delta_{c_0}, \delta_{c_1}, \dots, \delta_{c_{r-1}}]^T$ to M_C while $\Delta_U = 0$. The polluted adversarial measurements become $M_C^* = M_C + \Delta_C$, and $m_{c_i}^* = m_{c_i} + \delta_{c_i}$ ($0 \leq i \leq r-1$). Apparently, we have $\delta_i = \delta_{c_j}$ when $i = c_j$, $i \in C$, and $\delta_i = 0$ when $i \notin C$. Similarly, the crafted adversarial example $M^* = [m_0^*, m_1^*, \dots, m_{d-1}^*]^T = M + \Delta$ is fed into f_θ . We have $m_i^* = m_{c_j}^*$ when $i = c_j$, $i \in C$ and $m_i^* = m_i$ when $i \notin C$. All the notations are summarized in Table 1.

3.3.2 Mathematical Presentation. For **linear equality constraints**, such as the current measurements (Amperes) of the three meters in Figure 1, we suppose there are k constraints of the compromised

Table 1: List of Notations

Symbol	Description
f_θ	The trained model with hyperparameter θ
S	The vector of sensors
M	The vector of measurements of S
Δ	The perturbations vector added to M
M^*	The sum of Δ and M . The vector of compromised input
C	The vector of the indexes of compromised sensors or measurements
U	The vector of the indexes of uncompromised sensors or measurements
Y	The original class of the measurement M
Φ	The linear constraint matrix

measurements M_C that the attacker needs to meet, and the k constraints can be represented as follow:

$$\begin{cases} \phi_{0,0} \cdot m_{c_0} + \dots + \phi_{0,r-1} \cdot m_{c_{r-1}} = \phi_{0,r} \\ \phi_{1,0} \cdot m_{c_0} + \dots + \phi_{1,r-1} \cdot m_{c_{r-1}} = \phi_{1,r} \\ \dots \\ \phi_{k-1,0} \cdot m_{c_0} + \dots + \phi_{k-1,r-1} \cdot m_{c_{r-1}} = \phi_{k-1,r} \end{cases} \quad (1)$$

The above constraints can be represented as (2). We have $\Phi_{k \times r} = [\Phi_0, \Phi_1, \dots, \Phi_{k-1}]^T$, where $\Phi_i = [\phi_{i,0}, \phi_{i,1}, \dots, \phi_{i,r-1}]$ ($0 \leq i \leq k-1$), $\Phi_{i,j} = \phi_{i,j}$ ($0 \leq i \leq k-1, 0 \leq j \leq r-1$) and $\tilde{\Phi} = [\phi_{0,r}, \phi_{1,r}, \dots, \phi_{k-1,r}]^T$.

$$\Phi_{k \times r} M_C = \tilde{\Phi} \quad (2)$$

The attacker generates the perturbation vector Δ_C and adds it to M_C such that f_θ will predict the different output. Meanwhile, the crafted measurements $M_C^* = \Delta_C + M_C$ should also meet the constraints in (2) to avoid being noticed by the system operator or detected by the error checking mechanism.

Formally, the attacker who launches AML attacks needs to solve the following optimization problem:

$$\max_{\Delta_C} L(f_\theta(M^*), Y) \quad (3a)$$

$$s.t. \quad M_C^* = M_C + \Delta_C \quad (3b)$$

$$\Phi_{k \times r} M_C = \tilde{\Phi} \quad (3c)$$

$$\Phi_{k \times r} M_C^* = \tilde{\Phi} \quad (3d)$$

$$M^* = M + \Delta \quad (3e)$$

$$\Delta_U = 0 \quad (3f)$$

where L is a loss function, and Y is the original class label of the input vector M .

In addition, the **linear inequality constraints** among the compromised measurements can be represented as equation (4), and the constrained optimization problem to be solved is also similar to (3) but replacing (3c) with $\Phi_{k \times r} M_C \leq \tilde{\Phi}$ and (3d) with $\Phi_{k \times r} M_C^* \leq \tilde{\Phi}$ respectively.

$$\Phi_{k \times r} M_C \leq \tilde{\Phi} \quad (4)$$

4 DESIGN OF CONAML

The universal adversarial measurements algorithm is proposed in subsection 4.1 to solve the knowledge constraint of the attacker. Subsection 4.2 and subsection 4.4 analyze the properties of physical linear equality constraints and linear inequality constraints in AML respectively and present the adversarial algorithms. We set the maximum numbers of searching step in different algorithms to control the attack's time cost.

4.1 Universal Adversarial Measurements

Algorithm 1: Universal Adv-Measur Algorithm

```

1 Input:  $f_\theta, MU, M_C, \lambda, Y, MaxItera$ 
2 Output:  $M^*$ 
3 function uniAdvMeasur( $f_\theta, MU, M_C, \lambda, Y, MaxItera$ )
4   initialize  $\Delta = 0$ 
5   build set  $MUC = \{M_{C|U_0}, M_{C|U_1}, \dots, M_{C|U_N}\}$ 
6   set counter  $cycNum = 0$ 
7   while  $cycNum < MaxItera$  do
8     set  $flag$  to 0
9     for  $M_{C|U_i}$  in  $MUC$  do
10       $\Delta = \text{onePerturGenAlgorithm}(\Delta, M_{C|U_i})$ 
11      if  $\text{sampleEva}(f_\theta, Y, MUC, \Delta) < \lambda$  then
12        set  $flag$  to 1
13        break
14      end
15    end
16    if  $flag$  equals 1 then
17      break
18    end
19     $cycNum++$ 
20  end
21  return  $M^* = M + \Delta$ 
22 end

```

We first deal with the challenge of the attacker's limited knowledge on the uncompromised measurements M_U . This challenge is difficult to tackle since the complete measurement vector M is needed to obtain the gradient values in many AML algorithms [17, 26, 34, 35, 40]. In 2017, Moosavi-Dezfooli *et al.* proposed the universal adversarial perturbation scheme which generates image-agnostic adversarial perturbation [35]. The identical universal adversarial perturbation vector can cause different images to be misclassified by the state-of-the-art ML-based image classifiers with high probability. The basic philosophy of [35] is to iteratively and incrementally build a perturbation vector that can misclassify a set of images sampled from the whole dataset.

Inspired by their approach, we now present our universal adversarial measurements algorithm. We define an ordered set of N sampled uncompromised measurements $MU = \{M_{U_0}, M_{U_1}, \dots, M_{U_{N-1}}\}$, and use $M_{C|U_i}$ to denote the crafted measurement vector from M_C and the sampled uncompromised measurement vector M_{U_i} . Here, $M_{C|U_i}$ is a crafted measurement vector with $\|M_{C|U_i}\| = d$.

Algorithm 2: Sample Evaluation

```

1 Input:  $f_\theta, Y, MUC, \Delta$ 
2 Output: Classification Accuracy
3 function sampleEva( $f_\theta, Y, MUC, \Delta$ )
4   add perturbation  $\Delta$  to all vectors in  $MUC$ 
5   evaluate  $MUC$  with  $f_\theta$  and label  $Y$ 
6   return the classification accuracy of  $f_\theta(MUC)$ 
7 end

```

The uncompromised measurement vectors in MU can be randomly selected from the attacker's alternative dataset.

Algorithm 1 describes a high-level approach to generate adversarial perturbations regardless of uncompromised measurements. The algorithm first builds a set of crafted measurement vector MUC based on MU and M_C , and then starts an iteration over MUC . The iteration process is limited to $MaxItera$ times to control the maximum time cost. The purpose is to find a universal Δ that can cause a portion of the vectors in MUC misclassified by f_θ . The function **sampleEva** described in Algorithm 2 evaluates MUC and Y with the ML model f_θ and returns the classification accuracy. $\lambda \in (0, 1]$ is a constant chosen by the attacker to determine the attack's success rate in MUC according to Δ . During each searching iteration, algorithm 1 builds and maintains the perturbation Δ increasingly using an adversarial perturbation generation algorithms, as shown by Line 10 in Algorithm 1. We will propose our methods to handle this problem in the next subsections.

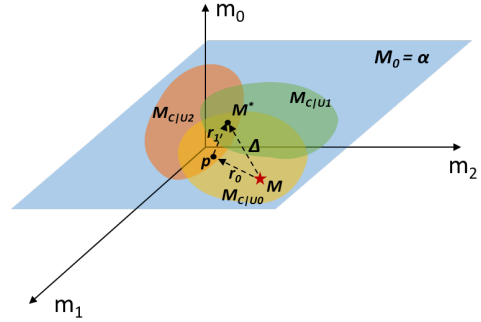


Figure 4: Iteration illustration.

Figure 4 presents a simple illustration of the iteration process in Algorithm 1. We assume there are three sensors' measurements $M = [m_0, m_1, m_2]$ in a CPS and only one sensor's measurement $m_0 = \alpha$ is compromised by the attacker. We set the sample number $N = 3$ and the yellow, green and orange shallow areas in the plane $M_0 = \alpha$ represent the possible adversarial examples of the crafted measurement vector $M_{C|U_0}$, $M_{C|U_1}$, and $M_{C|U_2}$, respectively, where U_i are randomly sampled measurements of uncompromised sensors (m_1 and m_2). The initial point M (red \star) iterates twice (r_0 and r_1) and finally reaches M^* with the universal perturbation vector Δ . Therefore, M^* is a valid adversarial example for all $M_{C|U_i}$ ($i \in \{0, 1, 2\}$).

Comparison of Methods: Our approach is different from [35] in several aspects. First, the approach proposed in [35] has identical adversarial perturbations for different ML inputs while our approach actually generates distinct perturbations for each M . Second, the approach in [35] builds universal perturbations regardless of the real-time ML inputs. However, as the attacker has already compromised a portion of measurements, it is more effective to take advantage of the obtained knowledge. In other words, our perturbations are ‘**universal**’ for M_U but ‘**distinct**’ for M . Finally, the intrinsic properties of CPSs require the attacker to generate a valid adversarial example within a sampling period while there is no enforced limitation of the iteration time in [35].

4.2 Linear Equality Constraints Analysis

As shown in [17] and [40], the fundamental philosophy of AML can be represented as (5).

$$M^* = M + \Delta = M + \epsilon \nabla_M L(f_\theta(M), Y) \quad (5)$$

However, directly following the gradient will not guarantee the adversarial examples meet the constraints in (2) and (4). With the constraints imposed by the physical system, the attacker is no longer able to freely add perturbation to original input using the raw gradient of the input vector. In this subsection, we will analyze how the linear equality constraints will affect the way to generate perturbation and use a simple example for illustration. The proofs of all the theorems and corollaries can be found in **Appendix B**.

Under the threat model proposed in Section 3.2, the constraint of (3c) is always met due to the properties of the physical systems. We then consider the constraint (3d).

THEOREM 4.1. *The sufficient and necessary condition to meet constraint (3d) is $\Phi_{k \times r} \Delta_C = 0$.*

From Theorem 4.1 we can also derive a very useful corollary, as shown below.

COROLLARY 4.2. *If $\Delta_{C_0}, \Delta_{C_1}, \dots, \Delta_{C_n}$ are valid perturbation vectors that follow the constraints, then we have $\Delta_{C'} = \sum_{i=0}^n a_i \cdot \Delta_{C_i}$ is also a valid perturbation for the constraint $\Phi_{k \times r}$.*

Theorem 4.1 indicates that the perturbation vector to be added to the original measurements must be a solution of the homogeneous linear equations $\Phi_{k \times r} X = 0$. However, is this condition always met?

THEOREM 4.3. *In practical scenarios, the attacker can always find a valid solution (perturbation) that meets the linear equality constraints imposed by the physical systems.*

We utilize a simplified example to illustrate how the constraints will affect the generation of perturbations, as shown in Figure 5. According to 5, measurement M should move a small step (perturbation) to the gradient direction (direction 1 in Figure 5) to increase the loss most rapidly. However, as shown by the contour lines in Figure 5, the measurement M is always forced to be on the straight line $y = 2 - 2x$ (2-dimension), which is the projection of the intersection of the two surfaces $z(x, y) = 2x^2 + 2y^2$ and $2x + y = 2$ (3-dimension). Accordingly, instead of following the raw gradient, M should move forward to direction 2 to increase the loss. Therefore, although at a relatively slow rate, it is still possible for the attacker to increase the loss under the constraints.

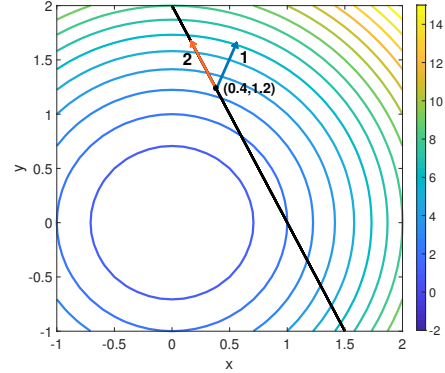


Figure 5: Linear equality constraint illustration. We consider a simple ML model f that only has two dimensions inputs (x, y) with a loss function $L(f_\theta(M), Y) = z(x, y) = 2x^2 + 2y^2$. Meanwhile, we suppose the input measurements x and y need to meet the linear constraints $2x + y = 2$ and the current measurement vector $M = (0.4, 1.2)$.

4.3 Adversarial Example Generation under Linear Equality Constraint

The common method of solving optimization problems using gradient descent under constraints is projected gradient descent (PGD). However, since neural networks are generally not considered as convex functions [8], PGD cannot be used to generate adversarial examples directly. We propose the design of a simple but effective search algorithm to generate the adversarial examples under physical linear equality constraints.

Algorithm 3: Best-Effort Search (Linear Equality)

```

1 Input:  $\Delta, f_\theta, C, M, step, size, \Phi, Y$ 
2 Output:  $v$ 
3 function genEqPer( $\Delta, f_\theta, C, M, step, size, \Phi, Y$ )
4   initialize  $v = \Delta$ 
5   initialize  $stepNum = 0$ 
6   while  $stepNum \leq step - 1$  do
7     if  $f'_\theta(M + v)$  doesn't equals  $Y$  then
8       | return  $v$ 
9     end
10     $r = \mathbf{eqOneStep}(f_\theta, C, M + v, size, \Phi, Y)$ 
11    update  $v = v + r$ 
12     $stepNum = stepNum + 1$ 
13  end
14  return  $v$ 
15 end

```

As discussed in subsection 4.2, the perturbation Δ_C needs to be a solution of $\Phi_{k \times r} X = 0$. We use $n = Rank(\Phi_{k \times r})$ to denote the rank of the matrix $\Phi_{k \times r}$, where $0 < n < r$. It is obvious that the solution set of homogeneous linear equation $\Phi_{k \times r} X = 0$ will have $r - n$ basic solution vectors. We use $I = [i_0, i_1, \dots, i_{r-n-1}]^T$ to denote the index

of independent variables in the solution set, $D = [d_0, d_1, \dots, d_{n-1}]^T$ to denote the index of corresponding dependent variables, and $B_{n \times (r-n)}$ to denote the linear dependency matrix of X_I and X_D . Clearly, we have $X_{D_{n \times 1}} = B_{n \times (r-n)} X_{I_{(r-n) \times 1}}$. For convenience, we will use $[I, D, B] = \mathbf{dependency}(\Phi_{k \times r})$ to describe the process of getting I, D, B from matrix $\Phi_{k \times r}$.

Algorithm 4: One Step Attack Constraint Δ_C

```

1 Input:  $f_\theta, C, M, size, \Phi, Y$ 
2 Output:  $r$ 
3 function eqOneStep( $f_\theta, C, M, size, \Phi_{k \times r}, Y$ )
4   calculate gradient vector  $G = \nabla_M L(f_\theta(M), Y)$ 
5   set all elements of  $G_U$  in  $G$  to zero
6   define  $G' = G_C$ 
7   obtain tuple  $[I, D, B] = \mathbf{dependency}(\Phi_{k \times r})$ 
8   update  $G'_D = B G'_I$  in  $G'$ 
9    $\epsilon = size / \mathbf{max}(\mathbf{abs}(G'))$ 
10  return  $r = \epsilon G$ 
11 end

```

As shown in Algorithm 3, the function **genEqPer** takes Δ as an input and outputs a valid perturbation v for M . Algorithm 3 keeps executing **eqOneStep** for multiple times defined by *step* to generate a valid v increasingly. Function **eqOneStep** performs a single-step attack for the input vector and returns a one-step perturbation r that matches the constraints defined by Φ , which is shown in Algorithm 4. Due to Corollary 4.2, Δ and v will also follow the constraints. To decrease the iteration time, similar to [34], the algorithm will return the crafted adversarial examples immediately as long as $f'_{\theta'}$ misclassifies the input measurement vector $M + v$, as shown by Line 7 in Algorithm 3.

The philosophy of function **eqOneStep** in algorithm 4 is very straightforward. From the constraint Matrix Φ , we can get the independent variables I , dependent variables D and the dependency matrix B between them. We will simply keep the gradient values of I and use them to compute the corresponding values of D (Line 8) so that the final output perturbation r will follow Φ . The constant *size* defines the largest modification of a specific measurement value in one iteration to control the search speed.

4.4 Adversarial Example Generation under Linear Inequality Constraint

Algorithm 5: Non-Constraint Perturbation.

```

1 Input:  $f'_{\theta'}, U, M, size, Y$ 
2 Output:  $r$ 
3 function freeStep( $f'_{\theta'}, U, M, size, Y$ )
4   calculate gradient vector  $G = \nabla_M L(f'_{\theta'}(M), Y)$ 
5   set elements in  $G_U$  to zero
6    $\epsilon = size / \mathbf{max}(\mathbf{abs}(G))$ 
7   return  $r = \epsilon G$ 
8 end

```

Linear inequality constraints are very common in real-world CPS applications, like the water flow constraints in Figure 3. Due to measurement noise, real-world systems usually tolerate distinctions between measurements and expectation values as long as the distinctions are smaller than predefined thresholds, which also brings inequality constraints to data. Meanwhile, a linear equality constraint can be represented by two linear inequality constraints. As shown in equation (4), linear inequality constraints define the valid measurement subspace whose boundary hyper-planes are defined by equation (2). In general, the search process under linear inequality constraints can be categorized into two situations. The first situation is when a point (measurement vector) is in the subspace and meets all constraints, while the second situation happens when the point reaches boundaries.

Algorithm 6: Best-Effort Search (Linear Inequality)

```

1 Input:  $\Delta, f'_{\theta'}, C, U, M, step, size, \Phi, \tilde{\Phi}, Y$ 
2 Output:  $v$ 
3 function genIqPer( $\Delta, f'_{\theta'}, C, U, M, step, size, \Phi, \tilde{\Phi}, Y$ )
4   initialize  $pioneer = \Delta, valid = pioneer$ 
5   initialize  $stepNum = 0$ 
6   initialize  $V$  as empty // violated constrain index
7   while  $stepNum \leq step - 1$  do
8     if  $f'_{\theta'}(M + valid)$  doesn't equals  $Y$  then
9       break
10    end
11     $chkRst = \mathbf{chkIq}(\Phi, \tilde{\Phi}, M + pioneer, C)$ 
12    if  $chkRst$  is empty then
13       $valid = pioneer$ 
14       $r = \mathbf{freeStep}(f'_{\theta'}, U, M + valid, size, Y)$ 
15       $pioneer = valid + r$ 
16      reset  $V$  to empty
17    else
18      extend  $V$  with  $chkRst$ 
19      define  $\Phi' = \Phi_V$  // real-time constraints
20       $r = \mathbf{eqOneStep}(f'_{\theta'}, C, M + valid, size, \Phi', Y)$ 
21       $pioneer = valid + r$ 
22    end
23     $stepNum = stepNum + 1$ 
24  end
25  return  $v = valid$ 
26 end

```

Due to the property of physical systems, the original point M will naturally meet all the constraints. As shown in Algorithm 6, to increase the loss, the original point will first try to move a step following the gradient direction through the function **freeStep** defined in Algorithm 5. Algorithm 5 is very similar to the FGM algorithm [40] but no perturbation is added to M_U , namely $r_U = 0$, which is similar to the saliency map function used in [38]. After that, the new point M' is checked with equation (4) to find if all inequality constraints are met. If all constraints were met, the moved step was valid and we can update $M = M'$. If M' violates some constraints

in Φ , we will take all the violated constraints and make a real-time constraint matrix Φ_V , where V is the index vector of violated constraints. We now convert the inequality constraint problem to the equality constraint problem with the new constraint matrix Φ_V and the original point M . M will then try to take a step using the **eqOneStep** function described in Algorithm 4 with the new constraint matrix Φ_V . Again, we check whether the new reached point meets all the constraints. If there are still violated constraints, we extend V with the new violated constraints. The search process repeats until reaching a valid M' that meets all the constraints. For simplicity, we will use $chkRst = \mathbf{chkIq}(\Phi, \tilde{\Phi}, M', C)$ to denote the checking process of a single search in one step movement, where $chkRst$ is the index vector of the violated constraints in the search.

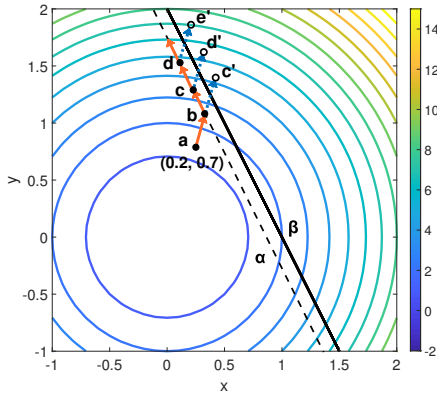


Figure 6: Best-Effort Search (linear inequality). We have the loss function $L(f_\theta(M), Y) = 2x^2 + 2y^2$ with inequality constraints $y \leq 2 - 2x$.

Similar to Figure 5, a simple example is shown in Figure 6. To increase the loss, the initial point a will take a small step following the gradient direction and reach point b . Since b meets the constraints, it is a valid point. After that, b will move a step following the gradient direction and reach point c' . However, point c' violates the constraint β and the movement is not valid. As we have point b is valid, we construct a linear equality constraint problem with constraint α which is parallel to β . With constraint α , point b will move a step to point c which is also a valid point. Point c then repeats the search process and increases the loss gradually. The real-time equality constraint is only used once. When a new valid point is reached, it empties the previous equality constraints and tries the gradient direction first.

5 EXPERIMENTAL EVALUATION

We evaluate our ConAML frameworks with two CPS study cases. The first study case is the ML-based false data injection (FDIA) detection in the power grids to show the impact of physical linear equality constraints, and the second case is the deep learning-based anomaly detection in the water treatment system to demonstrate linear inequality constraints.

Scenarios: For each CPS study, we consider four attack scenarios regarded to the different knowledge constraints, as summarized

in Table 2. The black-box scenario is the knowledge constraint we presented in our threat model in subsection 3.2 and is the most practical scenario for the attacker. We consider different scenarios to show the impact of different constraints and to study the robustness of ML models in CPSs under different circumstances. For white-box and gray-box1 scenarios, the attacker won't execute Algorithm 1 since the measurements of M_U are given and there is no knowledge constraint. We note that attackers under all the scenarios in Table 2 can only modify M_C and should also follow the physical and time constraints.

Table 2: Attack Scenarios

Scenario	Constraint
white-box	know both M_C and M_U , has access to f_θ
gray-box1	know both M_C and M_U , no access to f_θ
gray-box2	only know M_C , has access to f_θ
black-box	only know M_C , no access to f_θ

Baselines: We set two evaluation baselines in addition to the above scenarios. The first baseline is a supreme white-box attacker who has full access to the CPS ML model without considering any constraints and utilizes a state-of-the-art AML algorithm [40] to generate adversarial examples. We compare the performance of the ConAML framework with the supreme attack to demonstrate the impact of the constraints. The second baseline the autoencoder generator proposed in [11] by Erba *et al.* in 2019. In [11], an autoencoder is trained with the normal CPS measurement data and is expected to learn the physical constraints of measurements. The adversarial examples are fed to the autoencoder to be transferred into examples that meet the physical constraints. However, [11] allows the attacker to know complete M (same as the gray-box1 scenario in Table 2), which is less practical compared with our threat model. Meanwhile, since the patterns are learned by neural networks (autoencoder), the transferred measurements may not meet the linear constraints strictly. In our experiment, we show that the generated examples from the autoencoder may still violate the physical constraints.

Metrics: The evaluation metrics of ConAML can be different according to the attack purpose and the CPS properties. We set three metrics to evaluate the attack performance in this study. The first metric is detection accuracy of the defender's model under attack and a lower detection accuracy indicates a better attack performance. The second metric is the magnitude of the noise injected to the legitimate measurement. The attacker needs the adversarial examples to bypass the detection while maintaining their malicious behavior. A small bad noise will violate the attack's original intention even it can bypass the detection. We select the L_2 -Norm of the valid noise vector as the second metric to compare the magnitude of the malicious injected data. Finally, as the attack needs to be finished within a sampling period of the CPS, we will compare the time cost of the adversarial example generation.

5.1 Case Study: State Estimation in Power Grids

5.1.1 Background: State Estimation and FDIA. State estimation is a backbone of various crucial applications in power system control

that has been enabled by large scale sensing and communication technologies, such as SCADA. It is used to estimate the state of each bus, such as voltage angles and magnitudes, in the power grid through analyzing other measurements. A DC model of state estimation can be represented as (6), where \mathbf{x} is the state, \mathbf{z} is the measurement, and $\mathbf{H}_{m \times n}$ is a matrix that determined by the topology, physical parameters and configurations of the power grid.

$$\mathbf{z} = \mathbf{H}\mathbf{x} + \mathbf{e} \quad (6)$$

Due to possible meter instability and cyber attacks, bad measurements \mathbf{e} may be introduced to \mathbf{z} . To solve this, the power system employs a residual-based detection scheme to remove the error measurements [33]. The residual-based detection involves non-linear computation (L_2 -Norm), however, research has shown that a false measurement vector following linear equality constraints can be used to pollute the normal measurements without being detected. In 2009, Liu *et al.* proposed the false data injection attack (FDIA) that can bypass the residual-based detection scheme and finally pollute the result of state estimation [31]. In particular, if the attacker knows \mathbf{H} , she/he could construct a faulty vector \mathbf{a} that meets the linear constraint $\mathbf{B}\mathbf{a} = \mathbf{0}$, where $\mathbf{B} = \mathbf{H}(\mathbf{H}^T\mathbf{H})^{-1}\mathbf{H}^T - \mathbf{I}$, and the crafted faulty measurements $\mathbf{z} + \mathbf{a}$ will not be detected by the system. A detailed introduction of state estimation, residual-based error detection, and FDIA can be found in **Appendix C.1**

Many detection and mitigation schemes to defend FDIA are proposed, including strategical measurement protection [4] and PMU-based protection [47]. In recent years, detection based on ML, especially neural networks, become popular in the literature [3, 21, 23, 36, 37, 44, 46]. The ML-based detection does not require extra hardware equipment and achieve the state-of-the-art detection performance. However, in this section, we will demonstrate that the attacker can construct an adversarial false measurement vector \mathbf{z}_{adv} that can bypass both the residual-based detection and the ML-based detection. The ML models in previous research are trained to distinguish normal measurement \mathbf{z} and poisoned measurement $\mathbf{z} + \mathbf{a}$. Our ConAML algorithms allow the attacker to generate an adversarial perturbation \mathbf{v} that meets the constraint $\mathbf{B}\mathbf{v} = \mathbf{0}$ for his/her original false measurement $\mathbf{z} + \mathbf{a}$ and obtain a new adversarial false measurement vector $\mathbf{z}_{adv} = \mathbf{z} + \mathbf{a} + \mathbf{v}$ that will be classified as normal measurements by the ML-based FDIA detection models. The matrix \mathbf{B} then acts as the constraint matrix Φ defined in equation (3). Meanwhile, \mathbf{z}_{adv} can naturally bypass the traditional residual-based detection approach since the total injected false vector $\mathbf{a} + \mathbf{v}$ meets the constraint $\mathbf{B}(\mathbf{a} + \mathbf{v}) = \mathbf{B}\mathbf{a} + \mathbf{B}\mathbf{v} = \mathbf{0}$. Our experiment in the next subsection will show that our ConAML algorithms can significantly decrease the detection accuracy of the ML-based detection schemes.

5.1.2 Experiment Design and Evaluation. We select the IEEE standard 10-machine 39-bus system as the power grid system as it is one of the benchmark systems in related research [30, 36]. The features used for ML model training are the power flow (Ampere) measurements of each branches. The system has 46 branches so that there there will be 46 features for the ML models.

The goal of the attacker is to implement a false-negative attack that makes \mathbf{z}_{adv} bypass the detections. We utilize the MATPOWER [50] library to derive the \mathbf{H} matrix and simulate related datasets. We simulate two training datasets for the system operator and the

attacker to train their ML models. Through tuning the parameters, the overall detection accuracy of the defender's model f_θ is 98.3% and the attacker's model f'_θ is 97.5%. After that, we assume there are 10, 13, and 15 measurements being compromised by the attacker and simulate the corresponding test datasets that only contain the false measurements. A more detailed description of the experiment can be found in **Appendix C.2**.

Table 3: Evaluation Result Summary

Attack	Case	Accu	L_2 -Norm	Time (ms)
Supreme	10	0%	4077.43	5.8
	13	0%	8403.84	12.9
	15	0%	7979.26	6.8
Erba [11]	10	0%	1049.52	5.7
	13	0%	1164.71	5.84
	15	0%	1578.87	5.94
white-box	10	0%	2527.8	42
	13	0%	4984.03	96.8
	15	0%	7029.26	52.9
gray-box1	10	21.1%	2404.76	34.2
	13	48.9%	5356.09	87.1
	15	30.0%	9133.15	7.96
gray-box2	10	0%	2247.21	400.25
	13	5.4%	4882.95	222.4
	15	8.1%	6610.6	126.9
black-box	10	14.4%	1843.2	131.9
	13	4.3%	4786.72	209.6
	15	28.1%	9079.02	163.3

Table 3 summarizes the detection performance of f_θ under different adversarial attacks generated by our ConAML algorithms. From the table, we can learn that the ConAML attacks can effectively decrease the detection accuracy of the ML models used for FDIA detection and inject considerable bad data to the state estimation, even under black-box scenario. The autoencoder generator methods [11] can transfer the adversarial examples to follow the manifolds of the normal measurements (0% detection accuracy). However, the size of the successful bad data is very smaller compared with the supreme attack and ConAML, which decreases the effect of the FDIA attack. In addition, we check the adversarial examples generated by [11] and find that most of the generated examples (over 90%) of still violate the physical constraints and will be removed by the residual-based detection in state estimation.

As shown in Figure 7, by comparing the evaluation results of different cases, we can learn that compromising more sensors cannot guarantee better performances in attack detection. This is due to the different physical constraints imposed by the system. However, with more compromised sensors, the attacker can usually obtain a larger size of the injected bad data.

In our experiments, the time cost of gray-box2 and black-box is much higher than other attack scenarios due to the universal adversarial measurements algorithm, as shown in Figure 8. However, the time cost is still efficient for many CPS applications in practice. For example, the sampling period of the traditional SCADA system used in power systems is 2 to 4 seconds. In practical scenarios,

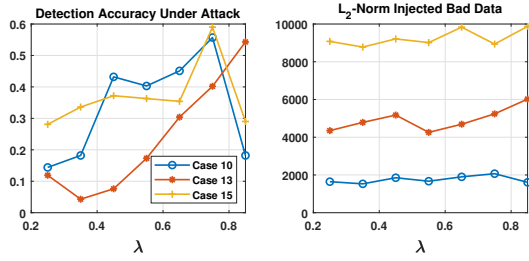


Figure 7: Performance of black-box attacks according to λ with $step = 40$, $size = 20$.

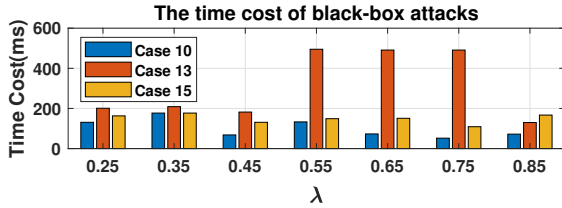


Figure 8: Time cost of black-box attacks according to λ with $step = 40$, $size = 20$.

the time cost also depends on the computational resource of the attacker. With the possible optimization and upgrade in software and hardware, the time cost can be further reduced.

5.2 Case Study: Water Treatment System

5.2.1 Background: SWaT Dataset. In this section, we study the linear inequality physical constraints based on the Secure Water Treatment (SWaT) proposed in [16]. SWaT is a scaled-down system but with fully operational water treatment functions. The testbed has six main processes and consists of cyber control (PLCs) and physical components of the water treatment facility. The SWaT dataset, generated by the SWaT testbed, is a public dataset to investigate the cyber attacks on CPSs. The raw dataset has 946,722 samples with each sample comprised of 51 attributes, including the measurements of 25 sensors and the states of 26 actuators. Each sample in the dataset was labeled with normal or attack. [16] investigated four kinds of attacks based on the number of attack points and places. The detailed description of the SWaT dataset can be found in [16] and [27].

The SWaT dataset is an important resource to study anomaly detection in CPSs. Inoue *et al.* used unsupervised machine learning, including Long Short-Term Memory (LSTM) and SVM, to perform anomaly detection based on the SWaT dataset [22]. By comparison, Kravchik *et al.* employed Convolutional Neural Networks (CNN) and achieved a better false positive rate [24]. In 2019, [13] proposed a data-driven framework to derive invariant rules for anomaly detection for CPS and utilized SWaT to evaluate their approach. Other research related to the SWaT dataset can be found in [1, 7, 12].

As shown in Table 4, the SWaT dataset includes the measurements from five kinds of analog components (25 sensors in total) whose measurements are used as the input features in previous

Table 4: SWaT Analog Components

Symbol	Description	Unit
LIT	Level Indication Transmitter	mm
FIT	Flow Indication Transmitter	m^3/hr
AIT	Analyzer Indication Transmitter	uS/cm
PIT	Pressure Indication Transmitter	kPa
DPIT	Differential Pressure Ind Transmitter	kPa

anomaly detection ML models. Our experiments aims to demonstrate that the ML models used for anomaly detection are vulnerable to adversarial attacks. However, due to the physical properties of the SWaT testbed, the sensor’s measurements are not independent but with linear inequality constraints.

In our experiment, we consider the scenario that the attacker compromises the FIT components to inject bad adversarial water flow measurements. We examined the SWaT testbed structure and find out that there are apparent linear inequality constraints among the FIT measurements. The linear inequality constraints of the seven FIT measurements in the dataset are defined by the structure of the water pipelines and the placement of the sensors, as shown in equation 7, where ϵ_1 and ϵ_2 are two constants of the system’s noise tolerance. We checked the SWaT dataset and observed that all the normal examples in the dataset meet the constraints. We also contacted the managers of the SWaT testbed and verified our find.

$$FIT301 \leq FIT201 \quad (7a)$$

$$\|FIT401 - FIT501\| \leq \epsilon_1 \quad (7b)$$

$$\|(FIT502 + FIT503) - (FIT501 + FIT504)\| \leq \epsilon_2 \quad (7c)$$

In our experiment, we show that the attacker can construct adversarial FIT measurements that can bypass the ML anomaly detection proposed in previous research. Meanwhile, our adversarial measurements will also follow the same linear inequality constraints to avoid being noticed by the system operator.

5.2.2 Experimental Design and Evaluation. Similar to the power system study case, we generate two training datasets for the defender’s model f_θ and the attacker’s model f'_θ , respectively by poisoning the normal measurements with Gaussian noise. The ML models are trained to distinguish the normal measurement data and the poisoned measurements (anomaly). In our experiment, the overall classification accuracy of f_θ and f'_θ is 97.2% and 96.7% respectively. After that, we consider the scenarios that there were 2, 5, and 7 FIT measurements compromised by the attacker and generate the related test datasets. The goal of the attacker is to generate the adversarial FIT measurements with the constraints defined by equation (7) so that the poisoned measurements can be classified as ‘normal’ by f_θ . A more detailed introduction of the experiment design and implementation, including the specific compromised measurements and the corresponding constraint matrix Φ , can be found in **Appendix D.1** and **D.2**.

Table 5 summarizes the evaluation performances of different scenarios of ConAML attacks. From the table, we can learn that

Table 5: Evaluation Result Summary

Attack	Case	Accu	L_2 -Norm	Time (ms)
Supreme	2	0%	3.24	3.59
	5	0%	2.85	6.58
	7	0%	4.31	9.21
Erba [11]	2	85.3%	0.176	4.72
	5	86.1%	0.017	4.12
	7	88.4%	0.184	6.18
white-box	2	0%	0.741	21.7
	5	0%	0.75	24.8
	7	2%	1.05	71.5
gray-box1	2	0%	0.522	21.4
	5	53.2%	0.466	51.0
	7	89.3%	0.835	136.4
gray-box2	2	1.0%	0.52	42.9
	5	1.2%	0.627	94.2
	7	1.3%	0.841	256.1
black-box	2	1.3%	0.309	17.5
	5	2.3%	0.340	111.7
	7	1.14%	0.411	451.8

the ConAML framework can still effectively decrease the detection accuracy of the ML models, even for black-box attacks. Meanwhile, even the black-box attack achieves a better performance on both the detection accuracy and bad data size compared with the baseline [11]. The size of the injected bad data of the ConAML attacks is smaller than the supreme attacker. We explain that this is due to the stringent constraints between the FIT measurements. Similar to the power system study case, a larger number of compromised sensors cannot produce a better performance in bypassing the detection. The reason for this result is that more compromised sensors will also have more complex constraints between their measurements. Meanwhile, more constraints will increase the computation overhead of the best effort search algorithms since there will be a ‘larger’ constraint matrix.

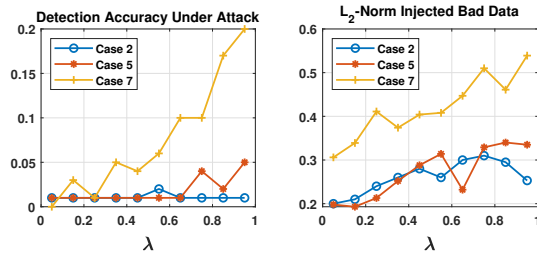


Figure 9: Performance of black-box attacks according to λ with $step = 50$, $size = 0.06$.

Figure 9 demonstrated the trend of the detection accuracy and injected bad data size according to λ . From the figure, we can learn that, with the λ increases, the probability of the adversarial examples being detected also increases. This matches the intuition that if an adversarial example can obtain higher successful attack probability with the sampling measurement set, its probability of

evading detection will also increase. Meanwhile, a smaller injected data size is expected to make the adversarial examples look more ‘normal’ to the detection model.

6 DISCUSSION AND FUTURE WORK

As we mentioned in Section 1, in this paper, we mainly investigate the linear constraints of input measurements in CPSs and neural network-based ML algorithms. In the future, research on ConAML of nonlinear constraints and other general ML algorithms, such as SVM, KNN will be proposed. We encourage related communities to present different CPSs that require special constraints.

As summarized in [48], defense mechanisms like adversarial re-training and adversarial detecting can increase the robustness of neural networks and are likely to mitigate ConAML attacks. However, most defenses in previous research target adversarial examples in computer vision tasks. In future work, we will study the state-of-the-art defense mechanisms in previous research and evaluate their performance with adversarial examples generated by ConAML. We will also investigate the defense mechanisms which take advantage of the properties of physical systems directly, such as the best deployment of sensors that will make the attackers’ constraint more stringent.

7 CONCLUSION

The potential vulnerability of ML applications in CPSs need to be concerned. In this paper, we investigate the input constraints of AML algorithms in CPSs. We analyze the difference of adversarial examples between CPS and computational applications, like computer vision, and give the formal threat model of AML in CPS. We propose the best-effort search algorithms to effectively generate the adversarial examples that meet the linear constraints. Finally, as proofs of concept, we study the vulnerabilities of ML models used in FDIA in power grids and anomaly detection in water treatment systems. The evaluation results show that even with the constraints imposed by the physical systems, our approach can still effectively generate the adversarial examples that will significantly decrease the detection accuracy of the defender’s ML models.

REFERENCES

- [1] Chudhry Mujeeb Ahmed, Jianying Zhou, and Aditya P Mathur. Noise matters: Using sensor and process noise fingerprint to detect stealthy cyber attacks and authenticate sensors in cps. In *Proceedings of the 34th Annual Computer Security Applications Conference*, pages 566–581. ACM, 2018.
- [2] T Athay, R Podmore, and S Virmani. A practical method for the direct analysis of transient stability. *IEEE Transactions on Power Apparatus and Systems*, (2):573–584, 1979.
- [3] Abdelrahman Ayad, Hany EZ Farag, Amr Youssef, and Ehab F El-Saadany. Detection of false data injection attacks in smart grids using recurrent neural networks. In *2018 IEEE Power & Energy Society Innovative Smart Grid Technologies Conference (ISGT)*, pages 1–5. IEEE, 2018.
- [4] Suzhi Bi and Ying Jun Zhang. Defending mechanisms against false-data injection attacks in the power system state estimation. In *2011 IEEE GLOBECOM Workshops (GC Wkshps)*, pages 1162–1167. IEEE, 2011.
- [5] Nicholas Carlini, Pratyush Mishra, Tavish Vaidya, Yuankai Zhang, Micah Sherr, Clay Shields, David Wagner, and Wenchao Zhou. Hidden voice commands. In *25th {USENIX} Security Symposium ({USENIX} Security 16)*, pages 513–530, 2016.
- [6] Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopadhyay. Adversarial attacks and defenses: A survey. *arXiv preprint arXiv:1810.00069*, 2018.
- [7] Yuqi Chen, Christopher M Poskitt, and Jun Sun. Learning from mutants: Using code mutation to learn and monitor invariants of a cyber-physical system. In *2018 IEEE Symposium on Security and Privacy (SP)*, pages 648–660. IEEE, 2018.

- [8] Anna Choromanska, Mikael Henaff, Michaël Mathieu, Gérard Ben Arous, and Yann LeCun. The loss surface of multilayer networks. *CoRR*, abs/1412.0233, 2014.
- [9] Joel Janek Dabrowski, Ashfaqur Rahman, Andrew George, Stuart Arnold, and John McCulloch. State space models for forecasting water quality variables: an application in aquaculture prawn farming. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 177–185. ACM, 2018.
- [10] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting adversarial attacks with momentum. In *Proceedings of the IEEE CVPR*, pages 9185–9193, 2018.
- [11] Alessandro Erba, Riccardo Taormina, Stefano Galelli, Marcello Pogliani, Michele Carminati, Stefano Zanero, and Nils Ole Tippenhauer. Real-time evasion attacks with physical constraints on deep learning-based anomaly detectors in industrial control systems. *arXiv preprint arXiv:1907.07487*, 2019.
- [12] Cheng Feng, Tingting Li, Zhanxing Zhu, and Deepthi Chana. A deep learning-based framework for conducting stealthy attacks in industrial control systems. *arXiv preprint arXiv:1709.06397*, 2017.
- [13] Cheng Feng, Venkata Reddy Pallethi, Aditya Mathur, and Deepthi Chana. A systematic framework to generate invariants for anomaly detection in industrial control systems. In *NDSS*, 2019.
- [14] Illinois Center for a Smarter Electric Grid. IEEE 39-Bus System. <https://icseg.iti.illinois.edu/ieee-39-bus-system/>. [Online; accessed 16-Aug-2020].
- [15] Amin Ghafouri, Yevgeniy Vorobeychik, and Xenofon Koutsoukos. Adversarial regression for detecting attacks in cyber-physical systems. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, pages 3769–3775, 7 2018.
- [16] Jonathan Goh, Sridhar Adepur, Khurum Nazir Junejo, and Aditya Mathur. A dataset to support research in the design of secure water treatment systems. In *International Conference on Critical Information Infrastructures Security*, pages 88–99. Springer, 2016.
- [17] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [18] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep recurrent neural networks. In *2013 IEEE international conference on acoustics, speech and signal processing*, pages 6645–6649. IEEE, 2013.
- [19] Kathrin Grosse, Nicolas Papernot, Praveen Manoharan, Michael Backes, and Patrick McDaniel. Adversarial examples for malware detection. In *European Symposium on Research in Computer Security*, pages 62–79. Springer, 2017.
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [21] Youbiao He, Gihan J Mendis, and Jin Wei. Real-time detection of false data injection attacks in smart grid: A deep learning-based intelligent mechanism. *IEEE Transactions on Smart Grid*, 8(5):2505–2516, 2017.
- [22] Jun Inoue, Yoriyuki Yamagata, Yuqi Chen, Christopher M Poskitt, and Jun Sun. Anomaly detection for a water treatment system using unsupervised machine learning. In *2017 IEEE International Conference on Data Mining Workshops (ICDMW)*, pages 1058–1065. IEEE, 2017.
- [23] JQ James, Yunhe Hou, and Victor OK Li. Online false data injection attack detection with wavelet transform and deep neural networks. *IEEE Transactions on Industrial Informatics*, 14(7):3271–3280, 2018.
- [24] Moshe Kravchik and Asaf Shabtai. Detecting cyber attacks in industrial control systems using convolutional neural networks. In *Proceedings of the 2018 Workshop on Cyber-Physical Systems Security and Privacy*, pages 72–83. ACM, 2018.
- [25] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *arXiv preprint arXiv:1607.02533*, 2016.
- [26] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial machine learning at scale. *arXiv preprint arXiv:1611.01236*, 2016.
- [27] ITrust Labs. Secure Water Treatment (SWaT) Dataset. https://itrust.sutd.edu.sg/itrust-labs_datasets, 2019. [Online; accessed 15-08-2019].
- [28] Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. Textbugger: Generating adversarial text against real-world applications. *arXiv preprint arXiv:1812.05271*, 2018.
- [29] Lei Li, Chieh-Jan Mike Liang, Jie Liu, Suman Nath, Andreas Terzis, and Christos Faloutsos. Thermocast: A cyber-physical forecasting model for datacenters. In *Proceedings of the 17th ACM SIGKDD, KDD '11*, pages 1370–1378. ACM, 2011.
- [30] Xuan Liu, Zhiyi Li, Xingdong Liu, and Zuyi Li. Masking transmission line outages via false data injection attacks. *IEEE Transactions on Information Forensics and Security*, 11(7):1592–1602, 2016.
- [31] Yao Liu, Peng Ning, and Michael K Reiter. False data injection attacks against state estimation in electric power grids. In *Proceedings of the 16th ACM conference on Computer and communications security*, pages 21–32, 2009.
- [32] Jiajun Lu, Hussein Sibai, Evan Fabry, and David Forsyth. No need to worry about adversarial examples in object detection in autonomous vehicles. *arXiv preprint arXiv:1707.03501*, 2017.
- [33] Alcir Monticelli. *State estimation in electric power systems: a generalized approach*. Springer Science & Business Media, 2012.
- [34] S. Moosavi-Dezfooli, A. Fawzi, and P. Frossard. Deepfool: A simple and accurate method to fool deep neural networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2574–2582, June 2016.
- [35] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1765–1773, 2017.
- [36] Xiangyu Niu, Jiangnan Li, Jinyuan Sun, and Kevin Tomsovic. Dynamic detection of false data injection attack in smart grid using deep learning. In *2019 IEEE Power & Energy Society Innovative Smart Grid Technologies Conference (ISGT)*, pages 1–6. IEEE, 2019.
- [37] Mete Ozay, Inaki Esnaola, Fatos Tunay Yarman Vural, Sanjeev R Kulkarni, and H Vincent Poor. Machine learning methods for attack detection in the smart grid. *IEEE transactions on neural networks and learning systems*, 27(8):1773–1786, 2015.
- [38] Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. In *2016 IEEE European Symposium on Security and Privacy (EuroS&P)*, pages 372–387. IEEE, 2016.
- [39] Nedin Rndic and Pavel Laskov. Practical evasion of a learning-based classifier: A case study. In *2014 IEEE symposium on security and privacy*, pages 197–211. IEEE, 2014.
- [40] Andras Rozsa, Ethan M Rudd, and Terrance E Boulton. Adversarial diversity and hard positive generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 25–32, 2016.
- [41] Mahmood Sharif, Sruti Bhagavatula, Lujo Bauer, and Michael K Reiter. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 1528–1540. ACM, 2016.
- [42] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [43] Yuchi Tian, Kexin Pei, Suman Jana, and Baishakhi Ray. Deeptest: Automated testing of deep-neural-network-driven autonomous cars. In *Proceedings of the 40th international conference on software engineering*, pages 303–314. ACM, 2018.
- [44] Chenguang Wang, Simon Tindemans, Kaikai Pan, and Peter Palensky. Detection of false data injection attacks using the autoencoder approach. *arXiv preprint arXiv:2003.02229*, 2020.
- [45] Allen J Wood, Bruce F Wollenberg, and Gerald B Sheblé. *Power generation, operation, and control*. John Wiley & Sons, 2013.
- [46] Jun Yan, Bo Tang, and Haibo He. Detection of false data attacks in smart grid with supervised learning. In *2016 International Joint Conference on Neural Networks (IJCNN)*, pages 1395–1402. IEEE, 2016.
- [47] Qingyu Yang, Dou An, Rui Min, Wei Yu, Xinyu Yang, and Wei Zhao. On optimal pmu placement-based defense against data integrity attacks in smart grid. *IEEE Transactions on Information Forensics and Security*, 12(7):1735–1750, 2017.
- [48] Xiaoyong Yuan, Pan He, Qile Zhu, and Xiaolin Li. Adversarial examples: Attacks and defenses for deep learning. *IEEE transactions on neural networks and learning systems*, 2019.
- [49] Zhenlong Yuan, Yongqiang Lu, Zhaoguo Wang, and Yibo Xue. Droid-sec: deep learning in android malware detection. In *ACM SIGCOMM Computer Communication Review*, volume 44, pages 371–372. ACM, 2014.
- [50] Ray Daniel Zimmerman, Carlos Edmundo Murillo-Sánchez, and Robert John Thomas. Matpower: Steady-state operations, planning, and analysis tools for power systems research and education. *IEEE Transactions on power systems*, 26(1):12–19, 2010.

A NON-LINEAR CONSTRAINTS

Many other ML applications in the CPS domains, for instance, load forecasting in power and water systems, traffic forecasting in transportation systems, may have nonlinear constraints. The nonlinear constraints can be very complex in various CPSs and cannot be covered in one study. In general, similar to linear constraints, the k nonlinear constraints of the compromised measurements can be represented as equation (8), where μ_i is a nonlinear function of M_C .

$$\begin{cases} \mu_0(m_{c_0}, m_{c_1}, \dots, m_{c_{r-1}}) = 0 \\ \mu_1(m_{c_0}, m_{c_1}, \dots, m_{c_{r-1}}) = 0 \\ \dots \\ \mu_{k-1}(m_{c_0}, m_{c_1}, \dots, m_{c_{r-1}}) = 0 \end{cases} \quad (8)$$

We now investigate a special case of the nonlinear constraints. If there exists a subset of the compromised measurements, in which

each measurement can be represented as an explicit function of the measurements in the complement set, the attacker will also be able to generate the perturbation accordingly. We use $P = [p_0, p_1, \dots, p_{n-1}]$ to denote the index vector of the former measurement set, and use $Q = [q_0, q_1, \dots, q_{r-n-1}]$ to denote the index vector of the complement set. We can then represent (8) as (9), where $\Xi = [\xi_0, \xi_1, \dots, \xi_{n-1}]$ is a vector of explicit functions.

$$\begin{cases} m_{p_0} = \xi_0(m_{q_0}, m_{q_1}, \dots, m_{q_{r-n-1}}) \\ m_{p_1} = \xi_1(m_{q_0}, m_{q_1}, \dots, m_{q_{r-n-1}}) \\ \dots \\ m_{p_{n-1}} = \xi_{n-1}(m_{q_0}, m_{q_1}, \dots, m_{q_{r-n-1}}) \end{cases} \quad (9)$$

Apparently, the roles of M_Q and M_P in (9) are similar to the M_I and M_D in linear constraints correspondingly. Instead of a linear matrix, the function set Ξ represents the dependency between M_P and M_Q . The nonlinear constraints make properties such as Theorem 1 infeasible. To meet the constraints, the attacker needs to find the perturbation Δ_Q first and obtain M_Q^* by adding it to M_Q . After that, the attacker can compute $M_P^* = \Xi(M_Q^*)$.

The above case of nonlinear constraints is special and may not be scalable to various practical applications. Although there are different types of nonlinear systems, they can be generalized using piece-wise linear constraints by setting proper ranges and break-points. We leave this as an open problem for future work.

B PROOFS

B.1 Theorem 4.1

PROOF. If we replace M_C^* in equation (3d) with equation (3b), we can get $\Phi_{k \times r} M_C^* = \Phi_{k \times r} (M_C + \Delta_C) = \Phi_{k \times r} M_C + \Phi_{k \times r} \Delta_C = \tilde{\Phi}$. From equation (3c) we can learn that $\Phi_{k \times r} M_C = \tilde{\Phi}$. Therefore, we have $\Phi_{k \times r} \Delta_C = 0$ and prove Theorem 4.1. \square

B.2 Corollary 4.2

PROOF. We have $\Phi_{k \times r} \Delta_{C'} = \Phi_{k \times r} \sum_{i=0}^n a_i \cdot \Delta_{C_i} = \sum_{i=0}^n a_i \cdot \Phi_{k \times r} \Delta_{C_i}$. Since Δ_{C_i} is a valid perturbation vector and $\Phi \Delta_{C_i} = 0$, we have $\Phi_{k \times r} \Delta_{C'} = 0$ and prove Corollary 4.2. \square

B.3 Theorem 4.3

PROOF. Due to the intrinsic property of the targeted system, equation (3c) is naturally met, which indicates that there is always a solution for the nonhomogeneous linear equations $\Phi_{k \times r} X = \tilde{\Phi}$. Accordingly, we have $\text{Rank}(\Phi_{k \times r}) \leq r$. Moreover, if $\text{Rank}(\Phi_{k \times r}) = r$, there will be one unique solution for equation (3c), which means the measurements of compromised sensors are constant. The constant measurements are contradictory to the purpose of deploying CPSs. In practical scenarios, M is changing over time, so that $\text{Rank}(\Phi_{k \times r}) < r$ and the homogeneous linear equation $\Phi_{k \times r} X = 0$ will have infinite solutions. Therefore, the attacker can always build a valid adversarial example that meets the constraints. \square

C POWER SYSTEM CASE STUDY

C.1 State Estimation and FDIA

We give the mathematical description of state estimation and how a false data injection attack (FDIA) can be launched. To be clear, we will employ the widely used notations in related research publications to denote the variables; the corresponding explanation will also be given to avoid confusion.

In general, the AC power flow measurement state estimation model can be represented as follow:

$$\mathbf{z} = \mathbf{h}(\mathbf{x}) + \mathbf{e} \quad (10)$$

where \mathbf{h} is a function of \mathbf{x} , \mathbf{x} is the state variables, \mathbf{z} is the measurements, and \mathbf{e} is the measurement errors. The task of state estimation is to find an estimated $\hat{\mathbf{x}}$ that best fits \mathbf{z} of (10). In practical application, a DC measurement model is also used to decrease the process time and (10) can then be represented as follow:

$$\mathbf{z} = \mathbf{H}\mathbf{x} + \mathbf{e} \quad (11)$$

where $\mathbf{H}_{m \times n}$ is a matrix that determined by the topology, physical parameters and configurations of the power grid.

Typically, if a weighted least squares estimation scheme is used, the system state variable vector $\hat{\mathbf{x}}$ can be obtained through (12):

$$\hat{\mathbf{x}} = (\mathbf{H}^T \mathbf{W} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{W} \mathbf{z} \quad (12)$$

where \mathbf{W} is the covariance matrix of the variances of meter errors.

Due to possible meter instability and cyber attacks, bad measurements may be introduced to the measurement vector \mathbf{z} . To solve this, various bad measurement detection methods are proposed [33]. One commonly used detection approach is to calculate the measurement residual between the raw measurement \mathbf{z} and derived measurements $\mathbf{H}\hat{\mathbf{x}}$. If the L_2 -norm $\|\mathbf{z} - \mathbf{H}\hat{\mathbf{x}}\| > \tau$, where τ is a threshold selected according to the false alarm rate, the measurement \mathbf{z} will be considered as a bad measurement.

The above detection method contains non-linear computation (L_2 -Norm), however, research has shown that a false measurement vector follows linear equality constraints can be used to pollute the normal measurements without being detected. In 2009, Liu *et al.* proposed the false data injection attack (FDIA) that can bypass the detection scheme described above and pollute the result of state estimation [31]. FDIA assumes that the attacker knows the topology and configuration information \mathbf{H} of the power system. Let $\mathbf{z}_a = \mathbf{z} + \mathbf{a}$ denote the compromised measurement vector that is observed by the state estimation, where \mathbf{a} is the malicious data added by the attacker. Thereafter, let $\hat{\mathbf{x}}_{bad} = \hat{\mathbf{x}} + \mathbf{c}$ denote the polluted state that is estimated by \mathbf{z}_a , where \mathbf{c} represents the estimation error brought by the attack. Liu *et al.* demonstrated that, as long as the attacker builds the injection vector $\mathbf{a} = \mathbf{H}\mathbf{c}$, the polluted measurements \mathbf{z}_a will not be detected by the measurement residual scheme.

PROOF. If the original measurements \mathbf{z} can pass the detection, the residual $\|\mathbf{z} - \mathbf{H}\hat{\mathbf{x}}\| \leq \tau$. Through (13) from [31], we learn that the measurement residual will be the same when $\mathbf{a} = \mathbf{H}\mathbf{c}$. Therefore, the crafted measurements from the attacker will not be detected. \square

$$\|\mathbf{z}_a - \mathbf{H}\hat{\mathbf{x}}_{bad}\| = \|\mathbf{z} + \mathbf{a} - \mathbf{H}(\hat{\mathbf{x}} + \mathbf{c})\| \quad (13a)$$

$$= \|\mathbf{z} - \mathbf{H}\hat{\mathbf{x}} + (\mathbf{a} - \mathbf{H}\mathbf{c})\| \quad (13b)$$

$$= \|\mathbf{z} - \mathbf{H}\hat{\mathbf{x}}\| \leq \tau \quad (13c)$$

Besides, [31] also provided the approach to effectively find vector \mathbf{a} that will meet the attack requirement. Let $\mathbf{P} = \mathbf{H}(\mathbf{H}^T\mathbf{H})^{-1}\mathbf{H}^T$ and matrix $\mathbf{B} = \mathbf{P} - \mathbf{I}$. In order to have $\mathbf{a} = \mathbf{H}\mathbf{c}$, \mathbf{a} needs to be a solution of the homogeneous equation $\mathbf{B}\mathbf{x} = \mathbf{0}$, as shown in (14).

$$\mathbf{a} = \mathbf{H}\mathbf{c} \Leftrightarrow \mathbf{P}\mathbf{a} = \mathbf{P}\mathbf{H}\mathbf{c} \Leftrightarrow \mathbf{P}\mathbf{a} = \mathbf{H}\mathbf{c} \Leftrightarrow \mathbf{P}\mathbf{a} = \mathbf{a} \quad (14a)$$

$$\Leftrightarrow \mathbf{P}\mathbf{a} - \mathbf{a} = \mathbf{0} \Leftrightarrow (\mathbf{P} - \mathbf{I})\mathbf{a} = \mathbf{0} \quad (14b)$$

$$\Leftrightarrow \mathbf{B}\mathbf{a} = \mathbf{0} \quad (14c)$$

Another problem of generating \mathbf{a} is when will (14c) have a solution. Liu *et al.* prove that, suppose the attacker compromises k meters, as long as $k > m - n$, there always exists non-zero attack vector $\mathbf{a} = \mathbf{H}\mathbf{c}$. We refer the readers to [31] for the detailed proof.

C.2 Experiment Implementation

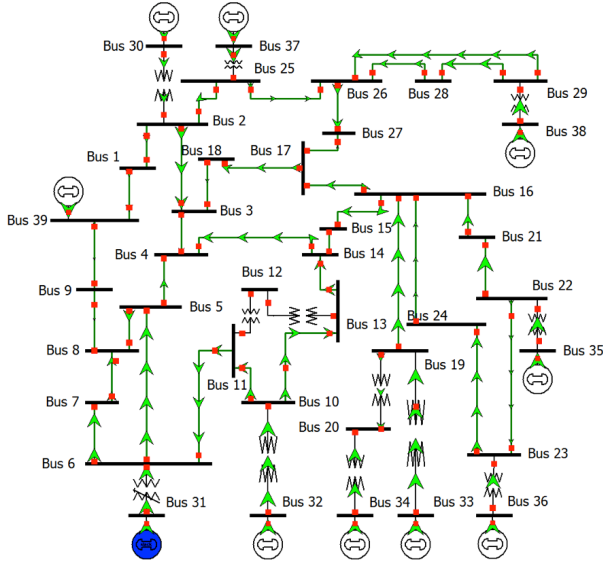


Figure 10: IEEE 39-Bus System [2] [14].

The structure of the IEEE 39-bus system is shown in Figure 10. We utilize the MATPOWER [50] library to derive the \mathbf{H} matrix of the system and simulate the power flow measurement data. We also implement the FDIA using MATLAB to generate false measurements. Both the power flow measurements and false measurements follow Gaussian distributions. We make two datasets for the defender and the attacker respectively. For each dataset, there are around 25,000 records with half records are polluted with FDIA. We label the normal measurements as 0 and false measurements as 1 and use one-hot encoding for the labels.

We investigate the scenarios that there are 10, 13, and 15 measurements being compromised by the attacker, with the randomly generated compromised index vector C and corresponding constraint matrix Φ (\mathbf{B}_C in (14)). We generate 1,000 false measurement vectors in each test datasets.

After that, we train two deep learning models based on the training datasets accordingly, with 75% records in the dataset used for training and 25% for testing. We use simple fully connected neural networks as the ML models and the model structures are shown in Table 6. Both the models are trained with a 0.0001 learning rate, 512 batch size, a mean squared error loss function, and a Stochastic Gradient Descent (SGD) optimizer. The deep learning models are implemented using Tensorflow and the Keras library and are trained on a Windows 10 machine with an Intel i7 CPU. The training process is around one minutes for each model.

Table 6: Model Structure - FDIA

Layer	f	f'
0	46 Input	46 Input
1	32 Dense ReLU	30 Dense ReLU
2	48 Dense ReLU	40 Dense ReLU
3	56 Dense ReLU	30 Dense ReLU
4	48 Dense ReLU	Dropout 0.25
5	32 Dense ReLU	20 Dense ReLU
6	Dropout 0.25	Dropout 0.25
7	16 Dense ReLU	2 Dense Softmax
8	Dropout 0.25	-
9	2 Dense Softmax	-

D WATER TREATMENT CASE STUDY

D.1 SWaT Measurement Constraints

We examined the user manual of the SWaT system and check the structure of the water pipelines. We found some FIT measurements in SWaT should always follow inequality constraints when the whole system is working steadily. Based on the component names described in [16], the constraints can be represented as (15), where ϵ_1 and ϵ_2 are the allowed measurement errors. We utilized the double value of the maximum difference of the corresponding measurements in the SWaT dataset to estimate ϵ_1 and ϵ_2 , and we had $\epsilon_1 = 0.0403$ and $\epsilon_2 = 0.153$.

$$\text{FIT301} \leq \text{FIT201} \quad (15a)$$

$$\|\text{FIT401} - \text{FIT501}\| \leq \epsilon_1 \quad (15b)$$

$$\|(\text{FIT502} + \text{FIT503}) - (\text{FIT501} + \text{FIT504})\| \leq \epsilon_2 \quad (15c)$$

Based on (4), we can represent (15) as follow. And M_C is the vector of measurements of FIT201, FIT301, FIT401, FIT501, FIT502, FIT503 and FIT504 accordingly.

$$\Phi_{5 \times 7} = \begin{bmatrix} -1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & -1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 & 1 & -1 \\ 0 & 0 & 0 & 1 & -1 & -1 & 1 \end{bmatrix} \quad \tilde{\Phi} = \begin{bmatrix} 0 \\ 0.0403 \\ 0.0403 \\ 0.153 \\ 0.153 \end{bmatrix}$$

We consider three scenarios that there are 2, 5, and 7 FIT measurements being compromised by the attacker, and the compromised sensors are {FIT201, FIT301}, {FIT401, FIT501, FIT502, FIT503, FIT504}, and all the seven FIT sensors respectively. The constraint matrix of each scenario can be derived from the corresponding rows of the $\Phi_{5 \times 7}$ matrix.

D.2 Experimental Implementation

In the Swat dataset, we extracted the normal records which were sampled when the whole system was working steadily. We also removed all the actuators' features. Here, we denote the extracted records as D_e . After that, we randomly picked out three test datasets from D_e as the with each test dataset contains 1000 records. We added Gaussian noise to the compromised measurements of records in all test datasets. We checked the polluted record every time when a noise vector was added to ensure all the records in test datasets meet the linear inequality constraints. Here, we denote the rest records of D_e as D_{train} which contains 120,093 records with each record having 25 features in our implementation. We randomly and equally split D_{train} into $D_{defender}$ and $D_{attacker}$ for the defender and attacker respectively and pollute half records with normally-distributed random noise in D_{train} and $D_{defender}$. The polluted records in D_{defend} and $D_{attacker}$ are labeled with 1 and the rest with 0. We allow the records in D_{train} and D_{defend} with label 1 to violate the constraints since the ML models are also expected to detect the obviously anomalous measurements.

We utilize D_{defend} and D_{attack} to train the ML models f_θ and f'_θ for the defender and attacker respectively. Again, 75% records in the both datasets were used for training the 25% records for testing. Similar to the FDIA experiment, we utilize fully connected neural networks and the structures are shown in Table 7. Through parameter tuning, model f_θ and f'_θ achieves 97.2% and 96.7% accuracy respectively.

Table 7: Model Structure - Water Treatment

Layer	f	f'
0	25 Input	25 Input
1	20 Dense ReLU	24 Dense ReLU
2	40 Dense ReLU	32 Dense ReLU
3	30 Dense ReLU	32 Dense ReLU
4	Dropout 0.25	16 Dense ReLU
5	20 Dense ReLU	2 Dense Softmax
6	Dropout 0.25	-
7	2 Dense Softmax	-