

Machine-learned Adversarial Attacks against Fault Prediction Systems in Smart Electrical Grids

Carmelo Ardito
carmelo.ardito@poliba.it
Polytechnic University of Bari
Bari, Italy

Yashar Deldjoo
yashar.deldjoo@poliba.it
Polytechnic University of Bari
Bari, Italy

Tommaso Di Noia
tommaso.dinoia@poliba.it
Polytechnic University of Bari
Bari, Italy

Eugenio Di Sciascio
eugenio.disciascio@poliba.it
Polytechnic University of Bari
Bari, Italy

Fatemeh Nazary*
fatemeh.nazary@poliba.it
Polytechnic University of Bari
Bari, Italy

Giovanni Servedio
g.servedio@studenti.poliba.it
Polytechnic University of Bari
Bari, Italy

ABSTRACT

In smart electrical grids, fault detection tasks may have a high impact on society due to their economic and critical implications. In the recent years, numerous smart grid applications, such as defect detection and load forecasting, have embraced data-driven methodologies. The purpose of this study is to investigate the challenges associated with the security of machine learning (ML) applications in the smart grid scenario. Indeed, the robustness and security of these data-driven algorithms have not been extensively studied in relation to all power grid applications. We demonstrate first that the deep neural network method used in the smart grid is susceptible to adversarial perturbation. Then, we highlight how studies on fault localization and type classification illustrate the weaknesses of present ML algorithms in smart grids to various adversarial attacks.

ACM Reference Format:

Carmelo Ardito, Yashar Deldjoo, Tommaso Di Noia, Eugenio Di Sciascio, Fatemeh Nazary, and Giovanni Servedio. 2023. Machine-learned Adversarial Attacks against Fault Prediction Systems in Smart Electrical Grids. In *AdvML '22: Workshop on Adversarial Learning Methods for Machine Learning and Data Mining, August 15, 2022, Washington DC, USA*. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION AND CONTEXT

Over the years, conventionally-operated electrical grids have undergone major revisions and upgrades in terms of dependability, robustness, and efficiency, thus moving to what we call today Smart Grids (SG). One of the most critical components of SGs is their application in fault detection, fault classification and routine examination of the underlying disruptions that trigger the failures. Power grid networks are inherently vulnerable to physical damages, and electrical faults can be caused by natural accidents such

as tree falling on a power line, a bird contact, lightning or aging of the equipment [19]. Power grid faults may lead to large-scale cascading effects, which might have a devastating impact on the economy and security of a country. As a result, rapid detection and classification of faults with a high degree of fidelity is a key service for the Electric Power Supply industry and the overall security of the critical energy infrastructure (CEI) [17].

This paper is focused towards classification of faults and their occurring area. Fault zone classification (FZC) aims to find the zone (or the exact location) in which the fault has occurred, while in fault type classification (FCT) the main objective is to determine the fault type class. Voltage sags are the main cause of faults, which can manifest as asymmetric phase-to-phase (LL), single-phase-to-ground (LG), or two-phase-to-ground (LLG) or symmetric three-phase-to-ground (LLLG or LLL) faults in both transmission and distribution systems [1, 22]. Previous literature has utilized a combination of tools and techniques from electrical engineering, signal processing, and artificial intelligence (AI) [9, 20, 22] to solve the above fault classification tasks. Among them, machine-learned (ML) models, notably those based on deep learning, have witnessed an increase in their acceptance in the current infrastructure of power systems, owing to the huge amounts of data spanning energy networks.

Unfortunately, notwithstanding their great performance, the intricacy of current (deep) inference methods may be their downfall. Adversarial attacks can take advantage of their vulnerabilities to compromise the confidentiality, integrity or availability of SGs (aka the CIA triad) [11, 25]. Adversarial attacks are operationalized via *adversarial examples* – subtle but non-random perturbations – designed to induce an ML model to produce erroneous outputs (e.g., to misclassify an input sample). As seen in Figure 1, an attacker can enter the SG system’s communication network in order to attack the failure prediction system used in supervisory control and data acquisition (SCADA) [8]. The purpose of the targeted adversarial attack is to induce the SCADA fault classification system’s machine learning (ML) model to misclassify an input sample as belonging to a known but incorrect class. Toward this aim, in the FZC scenario, the attacker selects an illegitimate target class label to prolong the rescue operation. By providing false positive signals, they can misdirect the rescue squad to places that do not need help.

Adversarial examples and attack on SG. Chen et al. [7] discuss the vulnerability of machine learning algorithms used in building load forecasting and power quality disturbance analysis against

*Corresponding author: Fatemeh Nazary (fatemeh.nazary@poliba.it). Authors are listed in alphabetical order.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
AdvML '22, August 15, 2022, Washington DC, USA

© 2023 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

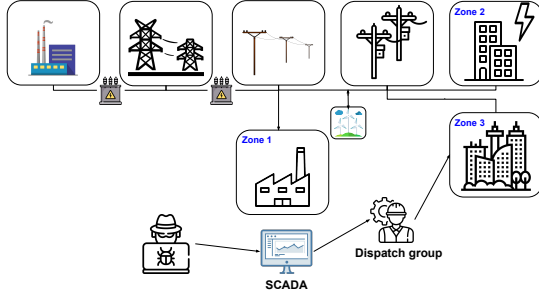


Figure 1: A hypothetical illustration of targeted adversarial attacks against fault zone prediction in smart grids. As a result of an adversarial attack on a fault location prediction system, dispatch recovery groups were dispatched to zone 3 by accident rather than zone 2, where they belonged.

specific adversarial attack. Farajzadeh-Zanjani et al. [12] developed a generative-adversarial system for partially labeled samples named semi-supervised generative adversarial learning. Adversarial attacks on convolutional neural network-based event causes have been presented in [15] for three specific power grid events: (1) line energization, (2) capacitor bank energization, and (3) fault prediction. The fast gradient sign technique (FGSM) [24] was used to introduce minor perturbations into voltage or current data for adversarial crafting. Additionally, the Jacobian-based Saliency Map Attack (JSMA) is used to compare the level of the FGSM’s adversary. Finally, adversarial training improves the CNN classifier’s performance against specific attacks. According to [23], voltage stability is assessed using adversarial instances generated using techniques such as FGSM, PGD, DeepFool, and Universal Adversarial Network (UAN), as well as Universal Adversarial Perturbation (UAP). Adversarial training is used to protect against these adversarial examples.

Our key contributions are summarized as follows:

- We investigate the impact of adversarial attacks against several key fault classification problems, and their combination on a widely used dataset based on the IEEE-13 test node feeder;
- We analyze adversarial attacks by examining multiple experimental situations with different adversarial goals;
- We show, via empirical experiments on a dataset collected from the IEEE-13 test node feeder, that adversarial attacks can degrade the quality of classifications significantly.

2 APPROACH

We have conducted adversarial attacks against two machine-learned fault classification task in smart electrical grids, which serve as the core attack target.

Adversarial task. Given a training dataset \mathcal{D} of n pairs $(x, y) \in \mathcal{X} \times \mathcal{Y}$, where x is the input sample, and y is its corresponding class label, the classification problem is formulated as finding a target function $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ that can predict the class label y surrounding the input sample x , where θ is the model parameter. The goal of the adversarial attacks is to find a non-random perturbation δ to produce an adversarial example $x^{adv} = x + \delta$ such that it can induce

an inaccurate detection (e.g., mis-classification). The methods by which δ is learned are referred to as *adversarial attacks*, and they can be either targeted or untargeted.

Definition 2.1 (Targeted adversarial attack). Given a trained classifier $f(x; \theta)$ and a test instance from the dataset $x_0 \in \mathcal{D}$ where $f(x_0; \theta) = y_0$, the goal of a targeted attack is to perturb x_0 with a small budget $\|\delta\| \leq \epsilon$ such that the perturbed sample would be mis-classified to the target label $y_T \neq y_0$, referred to as the *mis-classification label*. The problem can be represented using an unconstrained optimization problem formulation

$$\min_{\delta: \|\delta\| \leq \epsilon} \mathcal{L}(f(x_0 + \delta; \theta), y_T) \quad (1)$$

One can note that in this case, here the attacker aims to **minimize** the distance (loss) between the adversarial prediction $f(x_0 + \delta)$ and the mis-classification label y_T .

Definition 2.2 (Untargeted attack). The goal of the attacker in untargeted attack is to cause any mis-classification to *maximize* the loss between the adversarial prediction and the legitimate label y_0

$$\max_{\delta: \|\delta\| \leq \epsilon} \mathcal{L}(f(x_0 + \delta; \theta), y' \neq y_0) \quad (2)$$

as such, it is clear that the attacker’s objective in this scenario is to cause any mis-classification y' , regardless the of the specific type.

Fault Classification in Smart Grids. We consider different multi-class classification problems pertinent to fault prediction in smart grids with $K \geq 2$ classes in this paper, in which \mathcal{X} is the input space and $\mathcal{Y} = \{1, 2, \dots, K\}$ the output space. Our problem showcases two different target labels for the problems at hand: (i) fault location and (ii) fault type. Therefore, the main task is split into three sub-tasks:

- (1) Fault location classification (FLC): with $K = 4$ the task aims to classify a given signal into its originating zone as shown in Table 1.
- (2) Fault type classification (FTC): with $K = 11$ the task aims to classify a given signal into one of predefined fault types as shown in Table 1.
- (3) Joint location and type classification (FLC+FTC) $k = 44$ integrating the both fault class labels in the preceding cases;

where, (1) and (2) are explicitly contained in the dataset, while (3) is derived by combining each different possible combination of task (1) and task (2). Thus, we can state that the joint task is expected to be a more complex one compared to the former.

Adversary goal. The adversary is interested in mis-classifying smart-grid fault classification tasks in each of the three FLC, FTC, and joint sub-tasks through the use of two types of attacks: untargeted vs. targeted. In the latter situation, the purpose may be to produce more difficult-to-reach or difficult-to-resolve (mis-classification) labels in order to obstruct or delay the recovery task.

Adversary knowledge. Our assumption is a *white-box* setting where the attacker has full access to the feature extraction model parameters and input features that would be altered due to attacks. In a targeted attack scenario, the attacker can also obtain the class labels.

3 EXPERIMENTAL EVALUATION

We studied attacks on smart grids using data from the IEEE-13 test node feeder and explain the experimental setup below.

Table 1: The characteristic of the dataset used for training the machine-learned fault classification models in this work.

Item	Details
Fault type	phase to ground AG, BG, CG
	phase to phase AB, AC, BC
	phase to phase to ground ABG, ACG, BCG
	three phase ABC
	three phase to ground ABCG
Fault location	zone 1 branch 632-671
	zone 2 branch 632-633
	zone 3 branch 692-675
	zone 4 branch 671-680
Fault resistance	0.0010, 0.0273, 0.0535, 0.0798
	0.1061, 0.1323, 0.1586, 0.1848
	0.2111, 0.2374, 0.2636, 0.2899
	0.3162, 0.3424, 0.3687, 0.3949
	0.4212, 0.4475, 0.4737, 0.5, 1, 2

3.1 Datasets

For data collection and creating the training dataset for the fault classification in smart grids, similar to [1, 21, 22] we used short-circuit faults that were injected to IEEE-13 node test feeder using the MATLAB Simulink environment. The node feeder contained renewable energies such as wind turbine and photovoltaic system. We divided the network into four zones, adjacent to four load flow buses (numbered via 671, 633, 675, and 680, see [16]), and measured the three-phase voltage signals.

We applied 11 short circuit faults to four specified zone in the IEEE-13 network. These faults cover every conceivable short-circuit faults and are summarized in Table 1. To ensure having a sufficient number of samples in the training dataset, each fault was generated with 22 different fault resistance values [13, 21]. Our final training dataset contained $4 \text{ (zones)} \times 11 \text{ (faults)} \times 22 \text{ (resistance values)} \times 4 \text{ (measured locations)} = 3872$ samples. Note that we collected (measured) signals from 4 locations regardless of locations, and after feature extraction (see below) stacked them together to create a super-vector which was fed into the neural network ML model.

To inject faults, the entire simulation duration was carried out in the time interval $t = [0.0 - 0.022]$, with the network frequency 60Hz, sampling time 0.00001. Each fault with every resistance was applied at a certain start time $t = 0.01$ and revoked at a specified end time $t = 0.02$, hence $t_f = [0.01 - 0.02]$ represents the faulty duration and $t_h = [0 - 0.01]$ represents the healthy duration. For the signal type, in this work we only relied on (three-phase) voltage signals and kept investigation of other possible signals such as current for future investigation.

The time series signals were represented as discrete features retrieved from the time, frequency, and wavelet domains using temporal, Discrete Fourier transform (DFT), and Discrete wavelet transform (DWT) analysis, as previously explored [5, 18]. Afterwards, we extract from each domain, six features related to energy, maximum, as well as the 4-th moment of their probability distribution functions (PDFs) (e.g., mean, norm, skewness, kurtosis). The overall length of the feature vectors utilized in the learning model is 48, divided into 6 (time) + 6 (DFT) + 36 (DWT), where we employed 6 (coefficients) \times 6 (aggregation operations) for the DWT features, resulting in a 36-dimensional feature vector.

3.2 Adversarial Attacks

The performed attacks consist of the fast gradient sign method (FGSM), basic iterative method (BIM) [14], and Carlini and Wagner (C&W) [6]. FGSM is a white-box attack that employs the sign of the loss function's gradient to learn adversarial perturbations and BIM is the iterative version of the FGSM. Formally, in the untargeted scenario, FGSM aims to generate a perturbation that maximizes the training loss formulated as

$$\delta = \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \ell(f(\mathbf{x}; \theta), y)) \quad (3)$$

where ϵ (perturbation level) represents the attack strength and $\nabla_{\mathbf{x}}$ is the gradient of the loss function w.r.t. input sample \mathbf{x} , y is the legitimate label and $\text{sign}(\cdot)$ is the sign operator. A targeted FGSM attack is, instead, formulated as

$$\delta = -\epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \ell(f(\mathbf{x}; \theta), y_T)) \quad (4)$$

in which the goal of the attacker is maximize the conditional probability $p(y_T | \mathbf{x})$ for a given input \mathbf{x} .

The second category of adversarial attacks is Carlini and Wagner. It is a powerful attack model for finding adversarial perturbation under three various distance metrics (ℓ_0 , ℓ_2 , ℓ_∞). Its key insight is similar to L-BFGS [24] as it transforms the constrained optimization problem into an empirically chosen loss function to form an unconstrained optimization problem as

$$\min_{\delta} \left(\|\delta\|_p^p + c \cdot h(\mathbf{x} + \delta, y_T) \right) \quad (5)$$

where $h(\cdot)$ is the candidate loss function. \square

The C&W attack has been used with several norm-type constraints on perturbation ℓ_0 , ℓ_2 , ℓ_∞ among which the ℓ_2 and ℓ_∞ -bound constraint has been reported to be most effective [6].

4 EXPERIMENTS AND RESULTS

4.1 Explored Machine-Learnings Tasks

Model and training details. We trained a deep neural network, a Multi-layer Perceptron (MLP), for the three classification tasks specified in Section 2. The model is made of an input layer, two dense layers, and an output layer. The latter is the only layer that varies throughout the three tasks, as its number of neurons must correspond to the number of output classes in each task. The tasks require separate training phases, which all take place with the same settings, using 500 Epochs, Adam Optimizer, and fixed learning rate of $10e-3$ with a batch-size of 20. The hyper-parameters were obtained after fine-tuning.

Implementation of the attacks. We employed the IBM Adversarial Robustness Toolbox to perform the adversarial attacks due to its full compatibility with Keras and its wide offer of suitable attacks for a deep learning model. The performed attacks consist of FGSM, multi-step (BIM), and C&W attacks. These attacks were performed in both untargeted and targeted scenarios.

4.2 Results

Evaluation Questions. To obtain a better understanding of the effectiveness of the examined adversarial attacks against fault classification system in SGs, through the course of experiments, we intend to answer the following evaluation questions.

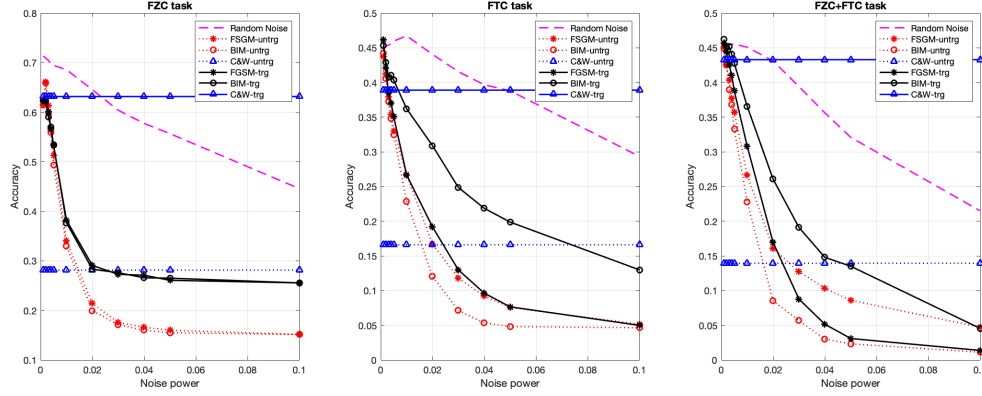


Figure 2: Three tasks under targeted and untargeted adversarial attacks. Classification accuracy for $FZC = 0.7134$, $FTC = 0.4569$, and $FZC + FTC = 0.4543$. Best results for C&W were obtained under ℓ_∞ for untargeted attacks and ℓ_2 for targeted attacks. Note that the starting point of noise power for all attacks and random noise is 0.001.

RQ 1: Against the three faults classification tasks in SGs presented in Section 2, how effective are adversarial perturbations generated by different adversarial attack methods (FGSM, BIM, and C&W) compared to random noise?

RQ 2: How does the performance of attacks change when we alternate between the **attack targets**?

Discussion. We begin our experimental study by addressing the above evaluation questions.

Answer to RQ 1. This research question verifies whether the application of adversarial attacks against fault classification system (FZC, FTC, and joint) has a sensible impact on the behavior of the ML models. As shown in Figure 2, all investigated adversarial attacks FGSM, BIM, and C&W have a much greater impact than random perturbation across three tasks and under different noise levels (ϵ), with the effect growing as the perturbation budget increases. Comparing the strength of the three adversarial attack models, BIM is the strongest in all tasks. For instance, in the case of (untargeted, FTC) with attack budget (noise level) equal to $\epsilon = 0.04$, BIM untargeted adversarial attack accuracy reaches 0.05, whilst FGSM and C&W reach 0.09 and 0.16, respectively, under the same condition. The effect of attack target (targeted vs. untargeted) is stronger on BIM and C&W than on FGSM. For example, for the FTC ($\epsilon = 0.04$), the classification accuracy is 0.21 vs. 0.05 (BIM-untargeted vs. BIM-targeted), while for FGSM the corresponding difference is only 0.1 vs. 0.09 (FGSM-untargeted vs. FGSM-targeted).

*In summary, the attacks' powers might be contrasted according to $BIM > C\&W > FGSM$ (the first being the strongest). The lone exception is **C&W-targeted**, which deviates from the trend and performs poorly, while **C&W-untargeted** performs well in all the explored scenarios.*

Answer to RQ 2. This research question verifies how much the performance of difference adversarial attacks varies across smart grid fault prediction tasks, and whether complexity of these tasks impacts the performances obtained.

We start this by assessing the absolute power of attacks across three tasks. At $\epsilon = 0.04$ the power of attacks FGSM-untarg, BIM-untargeted, C&W-untargeted, FGSM-targeted, BIM-targeted, C&W-targeted is equal to 0.166, 0.160, 0.281, 0.271, 0.265, and 0.631 respectively. Thus, w.r.t the base ML model (0.713), we may remark a relative degradation of 329%, 345%, 153%, 163%, 168%, and 13%. The equivalent relative degrading power of attacks for FTC task are 396%, 756%, 175%, 374%, 108%, 17% and for the joint FZC+FTC task include 339%, 1408%, 226%, 779%, 206%, 4.9%. Thus, the average degradation power for (untargeted, targeted) goals are, FZC=(275.6%, 114.6%), FTC=(442.3%, 166.3%), FZC+FTC=(657.6%, 329.9%). We might notice that both untargeted and targeted attack models work better (are stronger) as the task gets more complicated and this is true for both types of tasks.

In summary, the result of empirical evaluation shows that the complexity of the fault prediction tasks (in SGs) impacts the effectiveness of the explored adversarial attacks, meaning the attacks are better able to manipulate the decision outcomes according to $FZC+FTC > FTC > FZC$.

5 CONCLUSION

This work examines the security of fault classification systems in smart electrical grids powered by deep neural networks. Minor adversarial perturbations can reduce the quality of fault classification systems, highlighting the need for further studies to defend against adversarial training and detection methods (see [10]). Visual explanation of such adversarial threats [4] would constitute another interesting direction, which future work will investigate. Additionally, multi-party computation techniques, such as federated learning, could be used to develop privacy-preserving fault-prediction systems [2, 3], allowing separate zones to train models without exchanging data with a central server.

ACKNOWLEDGMENTS

This work has been partially funded by *e-distribuzione S.p.A* company, Italy, through a PhD scholarship granted to Fatemeh Nazary.

REFERENCES

- [1] Tamer S. Abdelgayed, Walid G. Morsi, and Tarlochan S. Sidhu. 2018. A New Harmony Search Approach for Optimal Wavelets Applied to Fault Classification. *IEEE Trans. Smart Grid* 9, 2 (2018), 521–529. <https://doi.org/10.1109/TSG.2016.2555141>
- [2] Vito Walter Anelli, Yashar Deldjoo, Tommaso Di Noia, and Antonio Ferrara. 2019. Towards effective device-aware federated learning. In *International Conference of the Italian Association for Artificial Intelligence*. Springer, 477–491.
- [3] Vito Walter Anelli, Yashar Deldjoo, Tommaso Di Noia, and Antonio Ferrara. 2020. Prioritized multi-criteria federated learning. *Intelligenza Artificiale* 14, 2 (2020), 183–200. <https://doi.org/10.3233/IA-200054>
- [4] Carmelo Ardito, Yashar Deldjoo, Tommaso Di Noia, Eugenio Di Sciascio, and Fatemeh Nazary. 2022. Visual inspection of fault type and zone prediction in electrical grids using interpretable spectrogram-based CNN modeling. *Expert Systems with Applications* 210 (2022), 118368.
- [5] Carmelo Ardito, Yashar Deldjoo, Eugenio Di Sciascio, and Fatemeh Nazary. 2021. Revisiting Security Threat on Smart Grids: Accurate and Interpretable Fault Location Prediction and Type Classification. In *Proceedings of the Italian Conference on Cybersecurity, ITASEC 2021, All Digital Event, April 7-9, 2021 (CEUR Workshop Proceedings, Vol. 2940)*, Alessandro Armando and Michele Colajanni (Eds.). CEUR-WS.org, 523–533. <http://ceur-ws.org/Vol-2940/paper44.pdf>
- [6] Nicholas Carlini and David A. Wagner. 2016. Defensive Distillation is Not Robust to Adversarial Examples. *CoRR* abs/1607.04311 (2016). arXiv:1607.04311 <http://arxiv.org/abs/1607.04311>
- [7] Yize Chen, Yushi Tan, and Deepjyoti Deka. 2018. Is Machine Learning in Power Systems Vulnerable?. In *2018 IEEE International Conference on Communications, Control, and Computing Technologies for Smart Grids, SmartGridComm 2018, Aalborg, Denmark, October 29-31, 2018*. IEEE, 1–6. <https://doi.org/10.1109/SmartGridComm.2018.8587547>
- [8] Lei Cui, Youyang Qu, Longxiang Gao, Gang Xie, and Shui Yu. 2020. Detecting false data attacks using machine learning techniques in smart grid: A survey. *J. Netw. Comput. Appl.* 170 (2020), 102808. <https://doi.org/10.1016/j.jnca.2020.102808>
- [9] Swagata Das, Sundaravaradan Navalpakkam Ananthan, and Surya Santoso. 2019. Estimating Zero-Sequence Line Impedance and Fault Resistance Using Relay Data. *IEEE Trans. Smart Grid* 10, 2 (2019), 1637–1645. <https://doi.org/10.1109/TSG.2017.2774179>
- [10] Yashar Deldjoo, Tommaso Di Noia, and Felice Antonio Merra. 2021. A survey on adversarial recommender systems: from attack/defense strategies to generative adversarial networks. *ACM Computing Surveys (CSUR)* 54, 2 (2021), 1–38.
- [11] Quang Do, Ben Martini, and Kim-Kwang Raymond Choo. 2016. A Data Exfiltration and Remote Exploitation Attack on Consumer 3D Printers. *IEEE Trans. Inf. Forensics Secur.* 11, 10 (2016), 2174–2186. <https://doi.org/10.1109/TIFS.2016.2578285>
- [12] Maryam Farajzadeh-Zanjani, Ehsan Hallaji, Roozbeh Razavi-Far, Mehrdad Saif, and Masood Parvania. 2021. Adversarial Semi-Supervised Learning for Diagnosing Faults and Attacks in Power Grids. *IEEE Trans. Smart Grid* 12, 4 (2021), 3468–3478. <https://doi.org/10.1109/TSG.2021.3061395>
- [13] Md Shakawat Hossain and Badrul H. Chowdhury. 2019. Data-Driven Fault Location Scheme for Advanced Distribution Management Systems. *IEEE Trans. Smart Grid* 10, 5 (2019), 5386–5396. <https://doi.org/10.1109/TSG.2018.2881195>
- [14] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. 2017. Adversarial examples in the physical world. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Workshop Track Proceedings*. <https://openreview.net/forum?id=HJGU3Rodl>
- [15] Iman Niaazari and Hanif Livani. 2020. Attack on Grid Event Cause Analysis: An Adversarial Machine Learning Approach. In *IEEE Power & Energy Society Innovative Smart Grid Technologies Conference, ISGT 2020, Washington, DC, USA, February 17-20, 2020*. IEEE, 1–5. <https://doi.org/10.1109/ISGT45199.2020.9087649>
- [16] Adeniyi Kehinde Onalapo, Kayode Timothy Akindeji, and Emmanuel Adetiba. 2019. Simulation experiments for faults location in smart distribution networks using ieee 13 node test feeder and artificial neural network. In *Journal of Physics: Conference Series*, Vol. 1378. IOP Publishing, 032021.
- [17] Ijeoma Onyeji, Morgan Bazilian, and Chris Bronk. 2014. Cyber security and critical energy infrastructure. *The Electricity Journal* 27, 2 (2014), 52–60.
- [18] Khaled A. Saleh, Ali Hooshyar, and Ehab F. El-Saadany. 2017. Hybrid Passive-Overcurrent Relay for Detection of Faults in Low-Voltage DC Grids. *IEEE Trans. Smart Grid* 8, 3 (2017), 1129–1138. <https://doi.org/10.1109/TSG.2015.2477482>
- [19] Enrico De Santis, Antonello Rizzi, and Alireza Sadeghian. 2018. A cluster-based dissimilarity learning approach for localized fault classification in Smart Grids. *Swarm Evol. Comput.* 39 (2018), 267–278. <https://doi.org/10.1016/j.swevo.2017.10.007>
- [20] Nikolaos Sapountzoglou, Jesus Lago, Bart De Schutter, and Bertrand Raison. 2020. A generalizable and sensor-independent deep learning method for fault detection and location in low-voltage distribution grids. *Applied Energy* 276 (2020), 115299.
- [21] Md Shafiullah and Mohammad A. Abido. 2018. S-Transform Based FFNN Approach for Distribution Grids Fault Detection and Classification. *IEEE Access* 6 (2018), 8080–8088. <https://doi.org/10.1109/ACCESS.2018.2809045>
- [22] Shenxing Shi, Beier Zhu, Sohrab Mirsaedi, and Xinzhou Dong. 2019. Fault Classification for Transmission Lines Based on Group Sparse Representation. *IEEE Trans. Smart Grid* 10, 4 (2019), 4673–4682. <https://doi.org/10.1109/TSG.2018.2866487>
- [23] Qun Song, Rui Tan, Chao Ren, and Yan Xu. 2021. Understanding Credibility of Adversarial Examples against Smart Grid: A Case Study for Voltage Stability Assessment. In *e-Energy '21: The Twelfth ACM International Conference on Future Energy Systems, Virtual Event, Torino, Italy, 28 June - 2 July, 2021*, Herman de Meer and Michela Meo (Eds.). ACM, 95–106. <https://doi.org/10.1145/3447555.3464859>
- [24] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. 2014. Intriguing properties of neural networks. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.).
- [25] Alireza Zarreh, HungDa Wan, Yooneun Lee, Can Saygin, and Rafid Al Janahi. 2019. Risk assessment for cyber security of manufacturing systems: A game theory approach. *Procedia Manufacturing* 38 (2019), 605–612.