

# Variational Auto-encoder-based Detection of Electricity Stealth Cyber-attacks in AMI Networks

Abdulrahman Takiddin, Muhammad Ismail, *Senior Member, IEEE*, Usman Zafar, Erchin Serpedin, *Fellow, IEEE*

**Abstract**—Current efforts to detect electricity theft cyber-attacks in advanced metering infrastructures (AMIs) are hindered by the lack of malicious electricity theft datasets. Therefore, anomaly detectors trained with the energy consumption profiles of honest customers appear as a plausible solution to overcome the lack of malicious datasets. Taking into account this constraint, this paper examines the performance of two structures of variational auto-encoders (VAEs); fully-connected (FC) VAE and long-short-term-memory (LSTM) VAE in detecting electricity thefts. The proposed structures are promising and exhibit an improvement of 11 – 15% in detection rate, 9 – 22% in false alarm rate, and 27 – 37% in the highest difference compared to existing state-of-the-art anomaly detectors that are shallow and static, such as single-class support vector machine (SVM) and auto-regressive integrated moving average (ARIMA) models.

**Index Terms**—electricity theft, auto-encoders, deep learning.

## I. INTRODUCTION

Electricity theft is a major problem for power companies not only because of the financial loss but also the grid overload and negative influence. To enable the energy consumption monitoring task in power grids, power companies are currently deploying advanced metering infrastructures (AMIs) with smart meters mounted in the customers' premises to track the energy consumption data. Regrettably, AMIs are subject to cyber electricity thefts. Malicious customers can hack AMIs and alter the integrity of energy consumption readings [1] – [4]. Thus, there is an acute need to develop algorithms to detect such cyber-attacks.

### A. Related Work and Limitations

Two classes of machine learning (ML) approached were proposed in the literature to identify electricity thefts. The first class is based on supervised learning and employs both benign and malicious energy consumption data for training. For example, the support vector machine (SVM)-based classifier in [5] exhibits a low detection accuracy of 72%. This approach suffers from the lack of malicious datasets to train the classifier. Moreover, it is not practical in case of attacks that the model is not trained to detect, such as zero-day attacks that

take place for the first time. Therefore, the second class, that is based on anomaly detection, is used. Anomaly detectors are trained using benign datasets to learn the honest customer's data consumption profile. Then, it identifies electricity theft patterns by measuring the deviation from the normal profile. For example, the auto-regressive integrated moving average (ARIMA)-based model presents a 77% detection rate [6].

Still, the existing anomaly detectors [1], [6], [7] suffer from three major weaknesses. First, the existing approaches rely mostly on shallow techniques with limited detection capabilities due to their inability to apprehend the complex patterns present in the energy consumption data. Second, most anomaly detectors do not make use of the temporal correlation in the energy consumption data. Third, the existing approaches exhibit a relatively low detection performance. Thus, new approaches are needed to capture the complex patterns and temporal correlations of the energy consumption readings to enhance the detection performance.

### B. Contributions

This paper proposes cyber-attack anomaly detectors based on variational auto-encoders (VAEs) and investigates their detection performance. The rationale behind the deep structure of the auto-encoders is to capture the complex patterns present in the energy consumption profiles. We consider sequence-to-sequence (seq2seq) structures based on long-short-term-memory (LSTM) recurrent neural networks (RNNs) to model the time-series nature of the energy consumption data. The contributions of this paper are as follows:

- VAE detectors are proposed to identify electricity cyber thefts. The advantage of using VAE lies in the latent variables, which are stochastic. Also, the probabilistic encoder of the VAE models the distribution of these latent variables. The variations in the latent space are captured by the VAE anomaly detector via the variance parameter.
- The performance of fully connected feed forward variational auto-encoders (FC-VAE) and LSTM-based RNN variational auto-encoders (LSTM-VAE) is investigated. FC-VAE presents a simple architecture with low computational complexity, while LSTM-VAE captures the temporal correlations in the energy consumption data [8].
- Sequential grid search hyper-parameter optimization [8] is adopted to optimize one hyper-parameter at a time to reduce the computational complexity and augment the overall detection performance.
- The proposed anomaly detectors are tested against six types of electricity theft cyber-attacks. The models' performance is compared to shallow architectures, namely,

A. Takiddin is with the ECEN Program, Texas A&M University at Qatar, Doha, Qatar, (email: abdulrahman.takiddin@qatar.tamu.edu).

M. Ismail is with the Department of Computer Science, Tennessee Tech University, Cookeville, TN, USA (email: mismail@tntech.edu).

U. Zafar is with Qatar Environment and Energy Research Institute, Hamad Bin Khalifa University, Doha, Qatar (email: uzafar@hbku.edu.qa).

E. Serpedin is with the ECEN Dept., Texas A&M University, College Station, TX, USA (e-mail: eserpedin@tamu.edu).

This publication was made possible by NPRP10-1223-160045 from the Qatar National Research Fund (a member of Qatar Foundation). The statements made herein are solely the responsibility of the authors.

single-class SVM and ARIMA. The developed VAE detector improves the detection rate (DR), false alarm (FA), and highest difference (HD) by 11–15%, 9–22%, and 27–37%, respectively.

The rest of this paper is structured as follows. Section II describes the benign and malicious load profiles and the data preparation step for training, validation, and testing. Section III explores the design and optimization of the anomaly detectors. Section IV depicts the experimental results. Section V concludes this paper's contributions.

## II. DATA PREPARATION

The data offered by the public Irish Smart Energy Trail [9] is adopted herein as the benign energy consumption data. The malicious data is created by employing the six cyber-attack models proposed in [1]. Each cyber-attack function models the electricity theft behavior of a malicious user.

### A. Benign Dataset

The Irish Smart Energy Trail dataset [9] is used to train and test the proposed electricity theft detector. This dataset consists of readings from 3,000 smart meters placed at residential units and that are recorded every half an hour for a 1.5-year period.

### B. Malicious Dataset

Let  $E_c(d, t)$  denote the energy consumption value for customer  $c$  at day  $d$  and time  $t$ . All these consumption values are regarded as entries of matrix  $\mathbf{E}_c$ . The energy consumption value reported by an honest customer's smart meter is denoted by  $R_c(d, t)$  ( $R_c(d, t) = E_c(d, t)$ ). Thus,  $\mathbf{E}_c = \mathbf{R}_c$  for honest customers. Malicious customers manipulate the integrity of the energy consumption readings and reduce their electricity bills such that  $R_c(d, t) \neq E_c(d, t)$ . We employ the false data injection approach [1] to build the malicious dataset. Next, we describe the set of cyber-attack functions, which can be classified into three main classes. The first class refers to *partial reduction attacks*. As a representative example of such a partial reduction attack, cyber-attack function  $f_1(E_c(d, t))$  decreases the actual energy consumption via a penalty constant factor  $\alpha$ , and the reported energy consumption  $R_c(d, t)$  is given by

$$R_c(d, t) = f_1(E_c(d, t)) = \alpha E_c(d, t). \quad (1)$$

As a more general example, cyber-attack function  $f_2(E_c(d, t))$  considers the dynamic penalty factor  $\beta(d, t)$ :

$$f_2(E_c(d, t)) = \beta(d, t) E_c(d, t). \quad (2)$$

The second class consists of *selective by-pass attacks*. In this class, malicious customers claim zero energy consumption at a given time window,  $[t_i(d), t_f(d)]$ , and report actual energy consumption data the rest of the time. Thus, it is modelled as:

$$f_3(E_c(d, t)) = \begin{cases} 0 & \forall t \in [t_i(d), t_f(d)] \\ E_c(d, t) & \forall t \notin [t_i(d), t_f(d)]. \end{cases} \quad (3)$$

The third class comprises the *price-based load control attacks*. These attacks are applicable to cases where the

electricity price varies during the day. In this set-up, one possible attack function may report a flat energy consumption value across the day:

$$f_4(E_c(d, t)) = \mathbb{E}[E_c(d)], \quad (4)$$

where  $\mathbb{E}[\cdot]$  denotes the expectation operator. To avoid reporting a constant value across the day, which can be easily detected, a dynamic fraction  $\beta(d, t)$  is employed:

$$f_5(E_c(d, t)) = \beta(d, t) \mathbb{E}[E_c(d)]. \quad (5)$$

Finally, another attack function that we consider reports high energy consumption values during periods of low electricity tariffs and vice versa:

$$f_6(E_c(d, t)) = E_c(d, T - t + 1). \quad (6)$$

Each of the above cyber-attack functions is applied to the customer energy consumption profile matrix  $\mathbf{E}_c$ . This operation leads to six malicious matrices per customer. The benign and malicious datasets are normalized to bring the values of all features to a common scale. The normalized dataset presents zero mean and unit variance. The normalized benign dataset  $\mathbf{B}$  is then divided into two disjoint subsets at the ratio 2:1. The first subset is used as training data  $\mathbf{X}_{\text{TR}}$ . The second subset is concatenated with the normalized malicious dataset  $\mathbf{M}$  to construct the test data. In this concatenation, each sample is associated with a label that takes value '0' if the sample is benign and value '1' if the sample is malicious. As much more malicious data is generated than benign data, and to avoid misleading performance results, we employ the adaptive synthetic sampling approach (ADASYN) [10] to balance the sets of benign and malicious data by over-sampling the minor (benign) class within the test set. Thus, we obtain test data  $\mathbf{X}_{\text{TST}}$  with label  $\mathbf{Y}_{\text{TST}}$ .

## III. DESIGN OF ELECTRICITY THEFT DETECTOR

Next, we focus on designing the electricity theft detectors. We analyze the adoption of two types of VAE architectures.

### A. Variational Auto-encoder Architecture

A VAE is a directed probabilistic graphical model whose posterior is approximated by a neural network [11]. The VAE assumes that a data point  $x$  is generated according to the unobserved continuous random variable  $k$ .  $k'$  represents a specific value that is generated using a prior distribution  $p(k)$ . Then,  $x'$ , which is the data instance is generated according to the conditional distribution  $p(x|k)$ . Since inferring the distribution  $p(k|x)$  is hard and the values of  $k$  are unknown, VAE determines the probability  $p(k|x)$  and  $p(x|k)$  through the network of encoders and decoders, respectively. The approximation of the true posterior  $p(k|x)$  is denoted by  $q(k|x)$  and the log-likelihood of data point  $x$  is

$$\log p(x) = D_{\text{KL}}(q(k|x)||p(k|x)) + \mathcal{L}(\Theta; x). \quad (7)$$

Notation  $D_{\text{KL}}$  designates the Kullback–Leibler (KL) divergence,  $\Theta$  denotes the model parameters, and  $\mathcal{L}(\Theta; x)$  is a variational lower bound on the log-likelihood expressed as

$$\mathcal{L}(\Theta; x) = -D_{\text{KL}}(q(k|x)||p(x)) + \mathbb{E}_{q(k|x)}[\log p_{\Theta}(x|k)]. \quad (8)$$

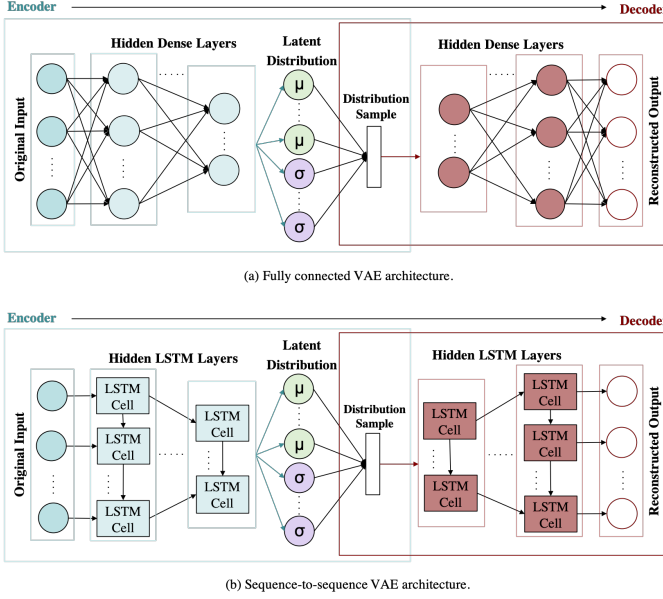


Fig. 1. Illustration of the VAE architecture.

The model parameters  $\Theta$  are learned by optimizing the lower bound  $\mathcal{L}(\Theta; x)$  ( $\mathcal{L}(\Theta; x) \leq \log p(k)$ ). If the latent variables are modeled as univariate Gaussian, then  $k = \mu + \sigma\mathcal{N}$  with  $\mathcal{N}$  denoting a normal distribution with zero mean and unit variance [15]. Therefore, the activations of the encoder forward-pass determine the mean ( $\mu$ ) and variance ( $\sigma^2$ ) of the Gaussian distribution. In VAE, the reconstruction probability is used as an anomaly score [12].

1) *Fully Connected VAE*: Figure 1a shows the structure of an FC-VAE. In VAE, the latent space is continuous and simplifies the random sampling for interpolation by outputting two vectors, mean vector  $\mu_x$  and variance vector  $\sigma_x^2$ . To generate samples of the energy consumption data, the decoder uses  $\mu_x$  and  $\sigma_x^2$ . Then, the reconstruction probability is calculated to detect electricity theft. The combined VAE loss function is expressed as follows

$$\hat{C} = \|x - \tilde{x}_\Theta\|^2 + D_{\text{KL}}(\mathcal{G}(\mu_x, \sigma_x^2) \parallel \mathcal{N}(0, 1)), \quad (9)$$

$\tilde{x}_\Theta$  represents the reconstructed output as function  $\Theta$ .  $\mathcal{G}(\mu_x, \sigma_x^2)$  and  $\mathcal{N}(0, 1)$  denotes the general and standard normal distribution, respectively.

The probability distribution's parameters,  $\mu_x$  and  $\sigma_x^2$ , are generated by the auto-encoder. The decoder then generates vectors that are used to construct an output via the decoder layers. The loss is then calculated using (9). After that, back propagation is used to update the parameters of the encoder and decoder. Algorithm 1 shows the calculation of the optimal parameters based on the iterative gradient descent.

2) *Sequence-to-Sequence VAE*: Originally, LSTM-VAEs were used as a way of data generation for language modeling applications [13]. Herein, we adopt an LSTM-VAE as an anomaly detector. Figure 1b shows the structure of an LSTM-VAE. The training algorithms of LSTM-VAE are supposed to generate  $\mu_x$  and  $\sigma_x^2$ , which are used to generate the samples that the decoder uses as shown in Algorithm 2.

### Algorithm 1: Training of FC-VAE

```

1 Input Data:  $X_{\text{TR}}$ 
2 Initialization: Weights  $W^l$  and biases  $b^l$  for all layers  $l$  and
  weights  $V_\mu$  and  $V_\sigma$  and biases  $b_\mu$  and  $b_\sigma$  for the latent layer
3 while not converged do
4   for each training sample  $x$  do
5     Feed Forward: Compute:
6     Encoder:
7     for all layers  $l = 1, \dots, L/2$  do
8        $z^l(x) = W^l a^{l-1}(x) + b^l$  and  $a^l(x) = \varphi(z^l(x))$ 
9     end
10    Generate  $\mu_x$  and  $\sigma_x^2$ :
11     $\mu_x = \varphi(V_\mu a^{L/2-1}(x)) + b_\mu^l$ 
12     $\sigma_x^2 = \varphi(V_\sigma a^{L/2-1}(x)) + b_\sigma^l$ 
13    Sample data  $\tilde{x}$  from  $\mathcal{G}(\mu_x, \sigma_x^2)$ 
14    Decoder:
15    for all layers  $l = L/2 + 1, \dots, L$  do
16       $z^l(\tilde{x}) = W^l a^{l-1}(\tilde{x}) + b^l$  and  $a^l(\tilde{x}) = \varphi(z^l(\tilde{x}))$ 
17    end
18    Back propagation: Compute:
19     $\nabla_{W^l_{(\cdot)}} \hat{C}$ ,  $\nabla_{V^l_{(\cdot)}} \hat{C}$  and  $\nabla_{b^l_{(\cdot)}} \hat{C}$ 
20  end
21  Weight and bias update:
22   $W^l_{(\cdot)} = W^l_{(\cdot)} - \frac{\eta}{K} \sum_x \nabla_{W^l_{(\cdot)}} \hat{C}$ 
23   $V^l_{(\cdot)} = V^l_{(\cdot)} - \frac{\eta}{K} \sum_x \nabla_{V^l_{(\cdot)}} \hat{C}$ 
24   $b^l_{(\cdot)} = b^l_{(\cdot)} - \frac{\eta}{K} \sum_x \nabla_{b^l_{(\cdot)}} \hat{C}$ 
25 end
26 Output: Optimal parameters  $W^l$ ,  $V_\mu$ ,  $V_\sigma$ ,  $b^l$ ,  $b_\mu$ , and  $b_\sigma$  for all
  layers

```

For both FC-VAE and LSTM-VAE, after the training is complete using  $X_{\text{TR}}$ , the test dataset  $X_{\text{TST}}$  is applied. Whenever the cost function that calculates the MSE between the original and reconstructed energy consumption profile is larger than a threshold, a malicious sample is labelled with  $y = '1'$ , else a benign sample is labelled with  $y = '0'$ .

### B. Performance Evaluation of the Detectors

TP, TN, FP, and FN define the true positives, true negatives, false positives, and false negatives, respectively. TP refers to a sample that is malicious and detected as malicious. TN indicates a benign sample that is detected as benign. FP means that the sample is benign but detected as malicious. FN represents the malicious sample that is detected as benign. To evaluate the performance of the developed detectors, we use (a) Detection Rate ( $\text{DR} = \text{TP}/(\text{TP} + \text{FN})$ ) that determines the number of malicious readings that were correctly detected as malicious by the detector. (b) False Alarm ( $\text{FA} = \text{FP}/(\text{TN} + \text{FP})$ ) refers to the number of benign samples that were incorrectly detected as malicious. (c) Highest Difference ( $\text{HD} = \text{DR} - \text{FA}$ ), which refers to the measured difference between DR and FA.

The calculated label  $Y_{\text{CAL}}$  is compared against  $Y_{\text{TST}}$  to produce a confusion matrix in order to calculate the performance evaluation metrics. To compute  $Y_{\text{CAL}}$ , a threshold is computed based on the median of the interquartile range (IQR) of the receiver operating characteristic (ROC) curve. If a score is less than the threshold value, it represents a benign sample. However, if a score is less than the threshold value, it represents a malicious sample.

---

**Algorithm 2: Training of LSTM-VAE**

---

```
1 Input Data:  $\mathbf{X}_{\text{TR}}$ 
2 Initialization: Weights  $\mathbf{U}_{(\cdot)}^l$ ,  $\mathbf{W}_{(\cdot)}^l$ ,  $\mathbf{V}_{(\cdot)}^l$ , and bias  $\mathbf{b}_{(\cdot)}^l \forall l$ 
3 while not converged do
4   for each training sample  $\mathbf{x}$  do
5     Feed Forward
6     Encoder:
7     for each hidden layer  $l = 1, \dots, L/2$  do
8       for each time step  $t$  do
9          $\mathbf{i}_{\text{E},t}^l = \varphi(\mathbf{W}_i^l \mathbf{x}_t^l + \mathbf{U}_i^l \mathbf{h}_{\text{E},t-1}^l + \mathbf{V}_i^l \mathbf{c}_{\text{E},t-1}^l + \mathbf{b}_i^l)$ ,
10         $\mathbf{f}_{\text{E},t}^l = \varphi(\mathbf{W}_f^l \mathbf{x}_t^l + \mathbf{U}_f^l \mathbf{h}_{\text{E},t-1}^l + \mathbf{V}_f^l \mathbf{c}_{\text{E},t-1}^l + \mathbf{b}_f^l)$ ,
11         $\mathbf{c}_{\text{E},t}^l = \mathbf{f}_{\text{E},t}^l \odot \mathbf{c}_{\text{E},t-1}^l + \mathbf{i}_{\text{E},t}^l \tanh(\mathbf{W}_c^l \mathbf{x}_t^l + \mathbf{U}_c^l \mathbf{h}_{\text{E},t-1}^l + \mathbf{V}_c^l \mathbf{c}_{\text{E},t-1}^l + \mathbf{b}_c^l)$ ,
12         $\mathbf{o}_{\text{E},t}^l = \varphi(\mathbf{W}_o^l \mathbf{x}_t^l + \mathbf{U}_o^l \mathbf{h}_{\text{E},t-1}^l + \mathbf{V}_o^l \mathbf{c}_{\text{E},t-1}^l + \mathbf{b}_o^l)$ ,
13         $\mathbf{h}_{\text{E},t}^l = \mathbf{o}_{\text{E},t}^l \tanh(\mathbf{c}_{\text{E},t}^l)$ ,
14      end
15       $\mathbf{h}'^l = \mathbf{h}_{\text{E},t}^l$ ,
16       $\mathbf{c}'^l = \mathbf{c}_{\text{E},t}^l$ .
17    end
18    Generate  $\mu_{\mathbf{x}}$  and  $\sigma_{\mathbf{x}}^2$ :
19     $\mu_{\mathbf{x}} = \varphi(\mathbf{V}_{\mu} \mathbf{a}^{l-1}(\mathbf{x})) + \mathbf{b}_{\mu}^l$ 
20     $\sigma_{\mathbf{x}}^2 = \varphi(\mathbf{V}_{\sigma} \mathbf{a}^{l-1}(\mathbf{x})) + \mathbf{b}_{\sigma}^l$ 
21    Sample data  $\tilde{\mathbf{x}}$  from  $\mathcal{G}(\mu_{\mathbf{x}}, \sigma_{\mathbf{x}}^2)$ 
22    Decoder:
23    The decoder hidden and cell states at initial time step are
    equal to  $\mathbf{h}'$  and  $\mathbf{c}'$ 
24    for each hidden layer  $l = L/2 + 1, \dots, L$  do
25      for each time step  $t$  do
26         $\mathbf{i}_{\text{D},t}^l = \varphi(\mathbf{W}_i^l \tilde{\mathbf{x}}_t^l + \mathbf{U}_i^l \mathbf{h}_{\text{D},t-1}^l + \mathbf{V}_i^l \mathbf{c}_{\text{D},t-1}^l + \mathbf{b}_i^l)$ ,
27         $\mathbf{f}_{\text{D},t}^l = \varphi(\mathbf{W}_f^l \tilde{\mathbf{x}}_t^l + \mathbf{U}_f^l \mathbf{h}_{\text{D},t-1}^l + \mathbf{V}_f^l \mathbf{c}_{\text{D},t-1}^l + \mathbf{b}_f^l)$ ,
28         $\mathbf{c}_{\text{D},t}^l = \mathbf{f}_{\text{D},t}^l \odot \mathbf{c}_{\text{D},t-1}^l + \mathbf{i}_{\text{D},t}^l \tanh(\mathbf{W}_c^l \tilde{\mathbf{x}}_t^l + \mathbf{U}_c^l \mathbf{h}_{\text{D},t-1}^l + \mathbf{V}_c^l \mathbf{c}_{\text{D},t-1}^l + \mathbf{b}_c^l)$ ,
29         $\mathbf{o}_{\text{D},t}^l = \varphi(\mathbf{W}_o^l \tilde{\mathbf{x}}_t^l + \mathbf{U}_o^l \mathbf{h}_{\text{D},t-1}^l + \mathbf{V}_o^l \mathbf{c}_{\text{D},t-1}^l + \mathbf{b}_o^l)$ ,
30         $\mathbf{h}_{\text{D},t}^l = \mathbf{o}_{\text{D},t}^l \tanh(\mathbf{c}_{\text{D},t}^l)$ ,
31      end
32    end
33    Back propagation: Compute  $\nabla_{\mathbf{W}_{(\cdot)}^l} C$ ,  $\nabla_{\mathbf{U}_{(\cdot)}^l} C$ ,
     $\nabla_{\mathbf{V}_{(\cdot)}^l} C$ , and  $\nabla_{\mathbf{b}_{(\cdot)}^l} C$ 
34  end
35  Weight and bias update:  $\mathbf{W}_{(\cdot)}^l = \mathbf{W}_{(\cdot)}^l - \frac{\eta}{K} \sum_{\mathbf{x}} \nabla_{\mathbf{W}_{(\cdot)}^l} C$ 
     $\mathbf{U}_{(\cdot)}^l = \mathbf{U}_{(\cdot)}^l - \frac{\eta}{K} \sum_{\mathbf{x}} \nabla_{\mathbf{U}_{(\cdot)}^l} C$ 
     $\mathbf{V}_{(\cdot)}^l = \mathbf{V}_{(\cdot)}^l - \frac{\eta}{K} \sum_{\mathbf{x}} \nabla_{\mathbf{V}_{(\cdot)}^l} C$ 
     $\mathbf{b}_{(\cdot)}^l = \mathbf{b}_{(\cdot)}^l - \frac{\eta}{K} \sum_{\mathbf{x}} \nabla_{\mathbf{b}_{(\cdot)}^l} C$ 
36 end
37 Output: Optimal  $\mathbf{U}_{(\cdot)}^l$ ,  $\mathbf{W}_{(\cdot)}^l$ ,  $\mathbf{V}_{(\cdot)}^l$ , and  $\mathbf{b}_{(\cdot)}^l \forall l$ .
```

---

### C. Hyper-parameter Optimization

Optimal choices of the detector hyper-parameters lead to improved detection performance. The hyper-parameters that we optimized are: the number of hidden layers (Dense or LSTM) ( $L$ ), which is the same for the encoder and decoder layers, the optimal number of neurons in those layers ( $N_l$ ), the optimizer ( $O$ ), the dropout rate ( $D$ ), and the hidden and output activation functions ( $A_{\text{H}}$  and  $A_{\text{O}}$ , respectively).

As shown in Algorithm 3, hyper-parameter optimization is carried out through four main sequential steps. Due to the large number of hyper-parameters that we aim to optimize, an exhaustive grid search presents a high computational complexity. Thus, we carried out a sequential grid search by optimizing one hyper-parameter at a time [8]. The motivation behind such an approach is to reduce the computational complexity and

improve the overall detection performance. For the selection of the hyper-parameters, we implement a cross-validation over  $\mathbf{X}_{\text{TR}}$  to decrease the chance of sub-optimality. Let  $P^*$  denote the hyper-parameter optimal setting that results in improved detection accuracy against the validation set. A given combination of hyper-parameters leads to a specific model (MD).

---

**Algorithm 3: Hyper-parameter Optimization**

---

```
1 Initialization: Optimizer = SGD, dropout rate = 0, hidden activation
   = Relu, output activation = Softmax
2 Output: A combination of optimized hyper-parameters
3 Input: Training set  $\mathbf{X}_{\text{TR}}$ 
4 for  $L \in \mathcal{L}$  do
5   for  $N_l \in \mathcal{N}$  do
6     Algorithms 1 and 2 are applied with  $L$  and  $N_l$  along with
7     other initial hyper-parameters ;
8     DR and FA are recorded;
9   end
10  The optimal  $L^*$  and  $N_l^*$  along with initial other hyper-parameters
   introduce model MD1
11 for  $O \in \mathcal{O}$  do
12   Algorithms 1 and 2 are applied with MD1's hyper-parameters
13   and  $o$ ;
14   DR and FA are recorded;
15 end
16  $L^*$ ,  $N_l^*$  and  $O^*$  along with initial other hyper-parameters introduce
   model MD2
17 for  $D \in \mathcal{D}$  do
18   Algorithms 1 and 2 are applied with MD2's hyper-parameters
19   and  $D$ ;
20   DR and FA are recorded;
21 end
22  $L^*$ ,  $N_l^*$ ,  $O^*$ , and  $D^*$  along with initial other hyper-parameters
   introduce model MD3
23 for  $A_{\text{H}} \in \mathcal{A}_{\text{H}}$  do
24   for  $A_{\text{O}} \in \mathcal{A}_{\text{O}}$  do
25     Algorithms 1 and 2 are applied with MD3's
26     hyper-parameters and  $A_{\text{H}}$  and  $A_{\text{O}}$ ;
27     DR and FA are recorded;
28   end
29 end
30  $L^*$ ,  $N_l^*$ ,  $O^*$ ,  $D^*$ ,  $A_{\text{H}}^*$ ,  $A_{\text{O}}^*$  are the optimal parameters.
```

---

## IV. EXPERIMENTAL RESULTS

### A. Threshold Value

After dividing the ROC curves into three quartiles and computing the IQR's median, the optimal threshold values for the FC-VAE and LSTM-VAE turned out to be 0.43 and 0.47, respectively.

### B. Hyper-parameter Optimization

Table I summarizes the optimized hyper-parameter values that are selected from the following sets: number of layers  $\mathcal{L} = \{2, 3, 4, 5\}$ , number of neurons  $\mathcal{N} = \{100, 200, 300, 400, 500\}$ , optimizer  $\mathcal{O} = \{\text{SGD}, \text{Adam}, \text{Adamax}, \text{Rmsprop}\}$ , dropout rate  $\mathcal{D} = \{0, 0.2, 0.4, 0.5\}$ , hidden activation functions  $\mathcal{A}_{\text{H}} = \{\text{Relu}, \text{Sigmoid}, \text{Linear}, \text{Tanh}\}$ , and output activation layer  $\mathcal{A}_{\text{O}} = \{\text{Softmax}, \text{Sigmoid}\}$ .

TABLE I  
OPTIMAL HYPER-PARAMETER VALUES

Hyper-parameter	FC-VAE	LSTM-VAE
$L^*$	8	4
$O^*$	Adam	SGD
$D^*$	0.4	0
$A_H^*$	Relu	Tanh
$A_O^*$	Softmax	Sigmoid

TABLE II  
PERFORMANCE EVALUATION

Model	DR	FA	HD
FC-VAE	88	11	77
LSTM-VAE	<b>91</b>	<b>7</b>	<b>84</b>
SVM	76	29	47
ARIMA	77	20	57

### C. Performance Evaluation

Table II summarizes the performance of the developed detectors. An improvement of 3% in DR, 4% in FA, and 7% in HD was observed when the LSTM-VAE model was used compared to the FC-VAE model. Such an improvement is due to the fact that the LSTM-based model captures better the time-series nature of the energy consumption data [14].

We also compare the performance of the developed deep auto-encoder-based anomaly detectors against the current state-of-the-art anomaly detectors, that are trained on benign data only, including (a) single-class SVM and (b) ARIMA-based anomaly detector that predicts future consumption with minimum prediction MSE. During testing, whenever the MSE is above a threshold, the detector announces a malicious sample. The SVM-based detector represents a static classifier that does not capture the time-series nature of the data. Although the ARIMA model captures the time-series nature of the data, it still represents a shallow architecture that does not capture well the complex patterns within the electricity readings. On the other hand, since the proposed VAEs represent a deep structure that captures the complex patterns and temporal correlations within the electricity consumption data, it outperforms the existing state-of-the-art anomaly detectors. As summarized in Table II, the developed VAE detector improves the DR, FA, and HD by 11 – 15%, 9 – 22%, and 27 – 37%, respectively.

## V. CONCLUSION

This paper proposed novel anomaly detectors for electricity theft detection based on variational auto-encoders. The developed anomaly detectors are trained only on benign energy consumption samples, a strategy that overcomes the limitation caused by the reduced number of malicious energy consumption profiles. This paper investigated whether deep architectures offer better detection performance compared to shallow detectors and whether recurrent LSTM-based architectures offer better detection performance compared to static fully connected feed forward-based detectors. Our study revealed a significant improvement when deep and recurrent anomaly detectors are employed compared to shallow and static structures. The best detection performance is achieved by LSTM-

VAE with 91% detection rate, 7% false alarm, and 84% highest difference offering an improvement of up to 15%, 22%, and 37% in detection rate, false alarm, and highest difference, respectively, compared to shallow detectors.

## REFERENCES

- [1] P. Jocar, N. Arianpoo, and V. Leung, "Electricity theft detection in AMI using customers' consumption patterns," *IEEE Trans. Smart Grid*, vol. 7, no. 1, Jan. 2016.
- [2] S. McIntyre, Termineter, Jan. 2018. [Online]. Available: <https://github.com/securestate/termineter/blob/master/README.md>
- [3] V. B. Krishna, C. A. Gunter, and W. H. Sanders, "Evaluating detectors on optimal attack vectors that enable electricity theft and der fraud," *IEEE J. Sel. Topics in Signal Processing*, vol. 12, no. 4, Aug. 2018.
- [4] Electric Sector Failure Scenarios and Impact Analyses Version 3.0, National Electric Sector Cybersecurity Organization Resource, Dec. 2015. [Online]. Available: <http://smartgrid.epri.com/doc/NESCOR-15.pdf>
- [5] J. Nagi, K. Yap, S. Tiong, S. Ahmed, and F. Nagi, "Improving SVM-based nontechnical loss detection in power utility using the fuzzy inference system," *IEEE Trans. Power Delivery*, vol. 26, no. 2, Apr. 2011.
- [6] V. Krishna, R. Iyer, and W. Sanders, "ARIMA-based modeling and validation of consumption readings in power grids," *10th International Conference on Critical Information Infrastructures Security (CRITIS 2015)*, pp. 199-210, May 2016.
- [7] J. Yeckle and B. Tang, "Detection of electricity theft in customer consumption using outlier detection algorithms," *1st International Conference on Data Intelligence and Security*, pp. 135-140, 2018.
- [8] I. Goodfellow, Y. Bengio, and A. Courville, "Deep Learning, MIT, 2016.
- [9] *Irish Social Science Data Archive*, Online: <http://www.ucd.ie/issda/data/commissionforenergyregulationcer/>.
- [10] H. He, Y. Bai, E. Garcia, and S. Li, "ADASYN: Adaptive synthetic sampling approach for imbalanced learning," *IEEE World Congress on Computational Intelligence*, pp. 1322-1328, June 2008.
- [11] D. Kingma and M. Welling, "Auto-encoding variational bayes," *2nd International Conf. Learning Representations (ICLR)*, pp. 1-14, 2014.
- [12] M. S. Kim, J. P. Yun, S. Lee and P. Park, "Unsupervised Anomaly detection of LM Guide Using Variational Autoencoder," *019 11th International Symposium on Advanced Topics in Electrical Engineering (ATEE)*, pp. 1-5, 2019.
- [13] S. Bowman, L. Vilnis, O. Vinyals, A. Dai, R. Jozefowicz, and S. Bengio, "Generating sentences from a continuous space," *20th Conf. Computational Natural Language Learning (CoNLL)*, pp. 10-21, 2016.
- [14] M. Nabil, M. Ismail, M. Mahmoud, M. Shahin, K. Qaraqe, and E. Serpedin, "Deep recurrent electricity theft detection in AMI networks with random tuning of hyper-parameters," *24th International Conference on Pattern Recognition (ICPR)*, Aug. 2018.
- [15] E. Serpedin, T. Chen, and R. Dinesh, *Mathematical Foundations for Signal Processing, Communications and Networks*, CRC Press/Francis & Taylor, 2012,