

Vulnerabilities of Power System Operations to Load Forecasting Data Injection Attacks

Yize Chen, Yushi Tan, Ling Zhang and Baosen Zhang

Abstract—We study the security threats of power system operation brought by a class of data injection attacks upon load forecasting algorithms. In particular, with minimal assumptions on the knowledge and ability of the attacker, we design attack data on input features for load forecasting algorithms in a black-box approach. System operators can be oblivious of such wrong load forecasts, which lead to uneconomical or even insecure decisions in commitment and dispatch. To our knowledge, this paper is the first attempt to bring up the security issues of load forecasting algorithms, and shows that accurate load forecasting algorithm is not necessarily robust to malicious attacks. More severely, attackers are able to design targeted attacks on system operations strategically with additional topology information. We demonstrate the impact of load forecasting attacks on two IEEE test cases. We show our attack strategy is able to cause load shedding with high probability under various settings in the 14-bus test case, and also demonstrate system-wide threats in the 118-bus test case with limited local attacks.

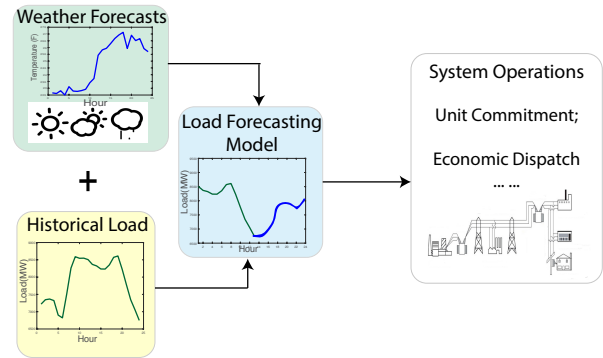
I. INTRODUCTION

Load forecasting plays an important role in the planning and operations of electric grids. As a cornerstone application for utilities and operators, it provides future load information which is utilized for various decision-making problems such as unit commitment, reserve management, economic dispatch and maintenance scheduling [1]. Consequently, the accuracy of forecasted loads directly impacts the cost and reliability of system operations [2].

Because of its fundamental importance, there are always strong incentives to improve short-term forecasting methods, especially under higher penetration of renewables. The driving factors of load variations are heterogeneous, including temperature, weather, temporal and seasonal effects (e.g., weekday vs. weekend) and other socioeconomic factors. Thus, load forecasting algorithms can be regarded as finding a nonlinear and complex mapping between the (potentially high dimensional) driving factors to the forecasted time series of load values. Over the past decades, a myriad of load forecasting algorithms have been proposed and adopted. See, for example, [1], [3] and the references within. Statistical and machine learning techniques, such as support vector regression [4], ARIMA [5] and neural networks [6] have been applied to short term load forecasting and implemented in practice. The recent advances in deep learning and data sciences opened the door to utilizing more input features and deeper model architectures to further improve load forecasting accuracy and provided some of the best performances to date [7]–[9]. In most of the load

forecasting studies, *forecast accuracy* has been regarded as the holy grail for researchers.

(a). Normal Operations



(b). Attack Construction

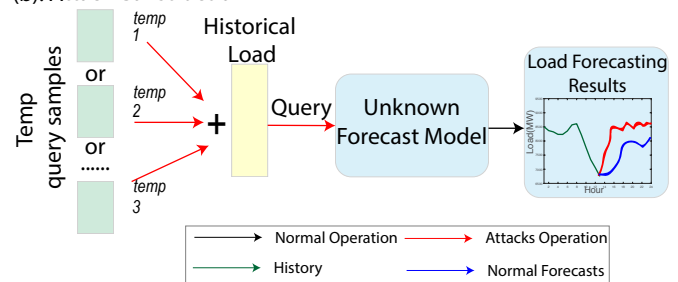


Fig. 1: The schematic of normal power system operations (a) and proposed black-box attacks on load forecasting (b). Without knowledge about the forecast model’s parameters, the attacker injects designed small, undetectable data perturbations into weather forecasts to induce abnormal system operations.

However, despite the critical role of forecasting algorithms and the long pursuit of *forecast accuracy*, the robustness and security issues have been overlooked to some extents. As the forecasting methods become more complex, they are also more susceptible to cybersecurity threats. In the previous research on cyber-security of power systems [10], [11], where state estimation [12], [13], communication [14], and electricity market [15], [16] threats and countermeasures are rigorously evaluated. However, the vulnerabilities of load forecasting algorithms are rarely discussed [17], [18], while this does not mean load forecasting is less vulnerable nor the consequences of attacks are less severe. For instance, forecasting models normally make use of weather forecasts inputs coming from external services/APIs, while such inputs can be exposed to adversarial modification and the model performance may be severely impacted by such malicious changes. Recently, there has been a hot debate on the security of machine learning models [19], and researchers found that small noises injected to the inputs can severely impact model performances [20].

Y. Chen, Y. Tan, L. Zhang and B. Zhang are with the Department of Electrical and Computer Engineering at the University of Washington, emails: {yizechen, ystan, lzhang18, zhangbao}@uw.edu. This work was supported by University of Washington Clean Energy Institute.

In this paper, we look into the security threats in general load forecasting algorithms. By taking the perspective of an attacker and developing attack strategies on load forecasting algorithms, we conduct damage analysis of the proposed attacks. We consider the scenario where attacker adversarially injects false data into the input features of forecasting algorithms, and examine to what extent such attacks could impact the performance of load forecasting models. Specifically, we investigate false data injection attacks on the temperature data. It is an important input to load forecasting algorithms and is mostly obtained from external services/APIs. Therefore it is easier for attackers to inject perturbations on temperature data than to attack the state estimation [13] or market clearing algorithms [15]. The potential damage of load forecasting attacks can be significant, leading to increases in system operation costs and maybe even more catastrophic events such as load shedding.

In Figure 1, we show the schematic of threats and proposed attacks to systems. Our work is different from most related work in two aspects: most of the studies in forecasting lack security and robustness considerations, while most of the studies in power system security evaluate attacks with certain level of knowledge about the targeted system or unconstrained capabilities. In contrast to previous attacks research that assume full knowledge of system configuration and strong capabilities of attackers on implementing attacks (see, e.g. [12], [13], [15]), we take a restrictive setting of both the attacker’s *knowledge* and *capabilities*. Under our setup, the attacker does not need to know any parameter of the targeted load forecasting algorithms, and could only inject constrained perturbations into input temperatures to avoid detection. We develop a simple data-driven attack strategy for finding the injected perturbations onto the original temperature data. Surprisingly, we find the proposed attacks significantly degrade the performance of a class of (accurate) load forecasting algorithms. With only few degrees of perturbations injected into input temperatures, the load forecasting algorithm’s output deviates drastically from original accurate forecasts.

To further illustrate security issues brought by load forecasting attacks, we embed the load forecasting algorithms into canonical power system operation case studies considering network constraints and power balance. We consider both cases when attacker possesses and does not possess additional information on system topology and parameters. For the former case where the attacker can strategically inject data perturbations under certain attack budgets, we design a greedy algorithm to compromise a subset of nodal load forecasts to cause targeted damages such as uneconomical generation, infeasible line flows and generator schedules, and load shedding. Simulations based on real-world load datasets on the IEEE 14-bus and 118-bus systems demonstrate the system operation vulnerabilities by only maliciously changing the temperature by a few degrees.

This study illustrates the need to look at other properties in addition to *forecast accuracy*, and the need for more comprehensive analysis when developing and applying load forecasting techniques. We demonstrate that accuracy may not mean robustness, and a wrong forecast of load potentially leads

to costly operation decisions or system damage. Specifically, we make the following contributions in this work:

- To the best of our knowledge, this work is the first to evaluate the security issues of load forecasting procedures in power system operations. Starting with the setup for load forecasting along with its role in power system operations (Sec. II), data vulnerabilities of current forecasting methods are formulated and discussed (Sec. III).
- Black-box attack algorithm *gradient estimation* is proposed to generate hard-to-detect, adversarial input data for load forecasting algorithms (Sec. III).
- We show that the strategically designed adversarial injections upon input features could target either increased system operating costs or load shedding. The resulting optimization formulation of the attack problem maybe of independent interest (Sec. IV).
- Case studies of power system operations on standard IEEE test cases using real-world load data reveal the prevalent vulnerabilities of current forecasting techniques and demonstrate potential damages on power system operations via proposed attacks (Sec. V).

Compared to our prior work [18] which works on single bus vulnerabilities analysis, we bring out both the load forecasting threats and the attack strategies on power networks. Extensive numerical simulations also verify such load forecasting vulnerabilities generally exist in power networks operations. We also make our code open source as a public package for evaluating load forecasting robustness and security¹. Due to the space limits of submission, we refer to the preprint for more details on attack implementations and thorough tests on various load forecasting algorithms [21].

II. PRELIMINARIES

In this section, we briefly describe the notations and setup for load forecasting algorithms, and illustrate how load forecasting serves as an important component of system operation in day-ahead commitment and real-time dispatch.

A. Load Forecasting

To set up and find parameters of the short-term load forecasting algorithm for a specific region, the system operator needs to collect a training dataset $\mathcal{D}_{tr} = \{(\mathbf{X}_{t-H}, \dots, \mathbf{X}_t); L_{t+k}\}_{t=1}^T$ based on available historical data. Here $L_{t+k} \in [0, 1]$ are scalars representing scaled nodal load values and H is the history horizon [1]. The model’s forecast horizon is denoted by k and ranges from one hour to one day in short-term forecasts. $\mathbf{X}_{t-i} \in [0, 1]^d$, $i = 0, \dots, H$ are scaled, d -dimensional input feature vectors. Feature vector \mathbf{X}_t includes historical records of load, weather forecasts including temperature, weather indicators (e.g., sunny, rainy or cloudy) and seasonal indicator variables such as weekdays/weekends and hour of the day [9]. We express it as $\mathbf{X}_t := \{L_t, \mathbf{X}_t^{temp}, \mathbf{X}_t^{index}\}$, where L_t is the load history records; \mathbf{X}_t^{temp} is the temperature value vector of current and neighboring regions, which could be acquired from either system historical records or weather

¹https://github.com/chennnnyize/load_forecasts_attack

forecast API; \mathbf{X}_t^{index} are a collection of indicators including weather characteristics, seasonal factors and time factors. In the task of load forecasting, one is interested to find a function parameterized by θ : $f_\theta(\cdot)$, which learns the mapping from $(\mathbf{X}_{t-H}, \dots, \mathbf{X}_t)$ to future loads \hat{L}_{t+k} . The mean absolute error (MAE) is widely used to measure the performance of f_θ , which is defined by the average L_1 norm of difference between forecasted loads \hat{L}_{t+k} and ground truth L_{t+k} . The schematic of load forecasting model is depicted in Figure 1(a).

We note that vulnerability analysis conducted by this paper is not constrained to certain forecasting algorithms. As long as the model output is sensitive with respect to input features, our proposed attack methods would be able to alter load patterns maliciously. For the discussion hereafter, we make use of an Recurrent Neural Networks (RNN) [6]–[9], which is a widely adopted forecast algorithm to model the temporal dependencies between feature inputs and forecasted values.

B. System Operations Model

For a N -bus power network with set of loads D and set of generators G , we consider the system operation setting consisting of a day-ahead planning stage using load forecasts and a real-time operational stage using actual load.

- 1) A unit commitment (UC) model considering reserve margins, startup and shutdown cost, minimum up/down time constraints and ramping constraints is used to create a commitment schedule $G_t, t = 1, \dots, T$ based on the day-ahead load forecasts $\hat{L}_t, t = 1, \dots, T$;
- 2) For each time t , the dispatch of the scheduled units \mathbf{p}_t and the actual dispatch costs $C_1(\mathbf{p}_t)$ are calculated according to a basic Economic Dispatch (ED) model [22] based on the actual load L_t and generation schedule G_t .

The actual daily operation costs are calculated via summing the 24-hour dispatch costs and the startup and shutdown costs. When ED based on the day-ahead commitment does not have a feasible solution, load shedding is used to maintain the balance between supply and demand. The shedded loads \mathbf{LS}_t also incur costs $C_2(\mathbf{LS}_t)$. Note that under perfect forecasts, the generation schedule G_t can minimize real-time dispatch costs.

III. ATTACK STRATEGIES

In this section, we first describe the objective and constraints for implementing load forecast attacks. We then illustrate how an attacker is able to design *white-box* attack with known load forecasting model parameters. Finally we describe under the *black-box* setting, how data injection attacks can be implemented even though the attacker only has limited query access to the load forecasting model.

A. Objective of Attacker

The attacker's goal is to distort the forecasted load as much as possible in a certain direction, e.g., to either increase or decrease forecasted values. Consider the task of training an accurate load forecasting models, where estimation of θ is

given by minimizing the L_1 -norm of the difference between model predictions and ground truth values:

$$\min_{\theta} \frac{1}{T} \sum_{t=1}^T \|f_\theta(\mathbf{X}_{t-H}, \dots, \mathbf{X}_t) - L_{t+k}\|_1 \quad (1)$$

where during training, ground truth of historical records on \mathbf{X}_t and L_{t+k} are used; during testing and real-world system implementations, we are using \mathbf{X}_t which are coming from weather forecast as input features. Once the model is learned, it can be applied in a rolling-horizon fashion.

In order to distort the output forecast values from the trained model based on (1), the attacker actually has two choices of inserting attacks: *to attack* \mathbf{X}_t or *to attack* θ . While trained model θ itself is often safely kept by the operators, system operator has to use external data such as weather forecasts \mathbf{X}_t^{temp} as input features for f_θ . This actually provides a backdoor for the attacker, whose goal is to inject perturbations into the weather forecasts coming from external services. By generating adversarial input data $\tilde{\mathbf{X}}_t^{temp}$ for $f_\theta(\cdot)$, model predictions are modified adversarially. We use $\gamma = \{-1, 1\}$ to denote the chosen attack direction by attackers. If $\gamma = 1$ ($\gamma = -1$), the attacker tries to find $\tilde{\mathbf{X}}$ to decrease (increase) the load forecasts values. Since load values are always positive, the attacker's goal is to find $\tilde{\mathbf{X}}$ that minimizes the value of $\gamma f_\theta(\tilde{\mathbf{X}}_{t-H}, \dots, \tilde{\mathbf{X}}_t)$.

B. Attacker's Knowledge

We consider two attack scenarios, *white box* and *black-box* attacks. In the *white-box* settings, the attacker is assumed to know exactly the model parameters θ . This is a strong assumption in the sense that load forecast model $f_\theta(\cdot)$ is fully exposed to the attacker. On the contrary, in the *black-box* setting, the attacker only knows which family of load forecasting model has been applied (e.g., NN or RNN), but is blind to the forecasting algorithms and has no knowledge of any parameters of f_θ . We assume the attacker only has *query* access to the load forecasting model². That is, the attacker could query the implemented load forecasting model by using different values of input features for a limited number of times, and then try to get insights on how f_θ works.

C. Attacker's Capability

From the attacker's perspective, it is necessary to construct attack injections while avoid being detected by the system operators' bad data detection algorithms. We consider several realistic constraints for attacker's capability: it could be upper bounded by the maximum number of perturbed entries in the input data, by the average deviations on all features, or by the largest deviation from the clean data. Mathematically, the attacker wants to keep $\|\tilde{\mathbf{X}}_t^{temp} - \mathbf{X}_t^{temp}\|_p$ bounded, where p can take different values such as 0, 1, ∞ to express certain norm constraints corresponding to different detection countermeasures.

²Such query access assumption is possible in many *forecast-as-a-Service* businesses, e.g., SAS energy forecasting and Itron forecasting.

In summary, we formulate the model of attackers as the following optimization problem:

$$\min_{\tilde{\mathbf{X}}_{t-H}^{temp}, \dots, \tilde{\mathbf{X}}_t^{temp}} \gamma f_\theta(\tilde{\mathbf{X}}_{t-H}, \dots, \tilde{\mathbf{X}}_t) \quad (2a)$$

$$s.t. \quad \|\mathbf{X}_{t-i}^{temp} - \tilde{\mathbf{X}}_{t-i}^{temp}\|_p \leq \epsilon, \quad i = 0, \dots, H \quad (2b)$$

Note that there is a parallel between the forecast problem (1) and attack problem (2), where the objective's optimization directions and optimization variables are exactly in the opposite directions: forecasting model works on model parameters to minimize forecast errors, while attacker works on model inputs to maximize the errors to targeted directions. However, due to lack of model knowledge in the black-box setting, it is a challenging task for attackers to find efficient attack input $\tilde{\mathbf{X}}^{temp}$ via (2). In the next subsection, we will show a black-box attack method generally working with attacker's knowledge coming from *query access* to the forecast algorithm.

D. Black-Box Attack

Under the case of *white-box* where model parameters are known to the attacker, it is possible to find the attack input via solving (2). For the convenience of notations, we omit the superscript on \mathbf{X} in some of the following paragraphs, and introduce the generalizable attack methods not only suitable for attacking temperature forecasts, but also suitable for injecting false data into other features.

Since most state-of-the-art load forecasting algorithms use complex models such as neural networks, the attacker's problem (2) is nonconvex and furthermore, there is no closed-form solution for $\tilde{\mathbf{X}}_{t-H}, \dots, \tilde{\mathbf{X}}_t$. Nevertheless, an attacker can still find some attack vectors iteratively by taking gradients with respect to each time step's temperature values. Even though this may not find the optimal solution to (2), because of the highly nonconvex nature of the forecasting model, a slight (suboptimal) perturbation of the input features would drastically change the forecast output.

Based on (2), we define a loss function \mathcal{L} with respect to each time step's feature $\tilde{\mathbf{X}}_{t-i}$, $i = 0, \dots, H$. Then the attacker iteratively takes gradients of \mathcal{L} to find the adversarial input $\tilde{\mathbf{X}}_{t-i}$. The constraints in (2b) is included in the loss function using a log-barrier:

$$\mathcal{L}(\tilde{\mathbf{X}}_{t-i}) = \gamma f_\theta(\tilde{\mathbf{X}}_{t-H}, \dots, \tilde{\mathbf{X}}_t) - \beta \log(\epsilon - \|\mathbf{X}_{t-i}^{temp} - \tilde{\mathbf{X}}_{t-i}^{temp}\|_p) \quad (3)$$

where β is the weight of the barrier term. Since there are a large number of parameters and input features in many load forecasting algorithms, it can be computationally expensive to compute the exact gradient values for each input feature. We follow a simpler method in [19] to only update the attack features based on the sign of the gradient at each iteration j :

$$\tilde{\mathbf{X}}_{t-i}^{(j+1)} = \tilde{\mathbf{X}}_{t-i}^{(j)} - \alpha \cdot \text{sign}(\nabla_{\tilde{\mathbf{X}}_{t-i}^{(j)}} (\mathcal{L}(\tilde{\mathbf{X}}_{t-i}^{(j)}))) \quad (4)$$

where α controls the step size for updating adversarial temperature values. The resulting adversarial temperature vector is obtained by applying (4) for several times.

Under the black-box setting where $\nabla_{\tilde{\mathbf{X}}_{t-i}^{(j)}} (\mathcal{L}(\tilde{\mathbf{X}}_{t-i}^{(j)}))$ is not able to compute, we assume attacker is able to query the load forecasting algorithm for a limited number of times, and it is still possible to construct adversarial temperature inputs by using queries to estimate the gradients. In Figure 1(b) we show the schematic on generating adversarial temperature instances via querying. For k -th dimension of the input feature at time stamp $t - i$, $\tilde{\mathbf{X}}_{k,t-i}^{j+1}$, the attacker needs to query the load forecasting system on each feature entry to calculate the two-sided estimation of the gradient of f_θ :

$$\nabla_{\tilde{\mathbf{X}}_{k,t-i}} f_\theta(\tilde{\mathbf{X}}) \approx \frac{f_\theta(\tilde{\mathbf{X}} + \delta \mathbf{e}_k) - f_\theta(\tilde{\mathbf{X}} - \delta \mathbf{e}_k)}{2\delta} \quad (5)$$

where \mathbf{e}_k is a d -dimensional vector with all zero except 1 at k -th component, and δ takes a small value for gradient estimation. Once the gradient is estimated for each dimension of temperature features, we can follow the same method of (4) to iteratively build the adversarial features using the estimated gradient vectors:

$$\tilde{\mathbf{X}}_{t-i}^{(j+1)} = \tilde{\mathbf{X}}_{t-i}^{(j)} - \alpha \gamma \cdot \text{sign}(\nabla_{\tilde{\mathbf{X}}_{t-i}^{(j)}} f_\theta(\tilde{\mathbf{X}}^{(j)})). \quad (6)$$

To satisfy norm constraints on the allowed perturbation of $\tilde{\mathbf{X}}$, the attacker projects the adversarial data back into the pre-defined norms after each iterative attack step. As shown in Figure 2, the output load forecasts deviate a lot from the ground truth values, while temperature perturbations are constrained within attacker's capability. In [23], techniques on reducing number of queries are discussed for attacking an image classifier, which may also further improve the query efficiency of load forecasting attacks.

IV. ATTACKS ON SYSTEM OPERATIONS

In this section, we illustrate the attacks on load forecasting input features could further threaten power system operations, and propose a realistic attack strategy for the attacker under attack budgets on number of compromised nodal forecasts.

A. Attack Objectives

We assume the attacker could only inject constrained attacks in the day-ahead planning stage. Under the day-ahead load forecasts \hat{L}_d^t for $|D|$ loads in the networks, the UC problem is to find the generation schedule and dispatch for $|G|$ generators while satisfying reserve and system operation constraints:

$$\min_{\mathbf{u}, \mathbf{p}} C(\mathbf{p}) + S(\mathbf{u}) \quad (7a)$$

$$s.t. \quad \sum_{g \in G} p_g^t = \sum_{d \in D} \hat{L}_d^t, \quad \forall t \in T \quad (7b)$$

$$u_g^t p_g^{\min} \leq p_g^t \leq u_g^t p_g^{\max}, \quad \forall g \in G, \quad \forall t \in T \quad (7c)$$

$$f_l^{\min} \leq f_l^t \leq f_l^{\max}, \quad \forall l \in F, \quad \forall t \in T \quad (7d)$$

$$\sum_{g \in \mathcal{G}(d)} p_g^t + \sum_{k \in \mathcal{K}(d)} f_k^t = \hat{L}_d^t, \quad \forall d \in D, \quad t \in T \quad (7e)$$

$$u_g^t - u_g^{t-1} = z_g^t - y_g^t, \quad \forall g \in G, \quad t \in T \quad (7f)$$

$$\sum_{\tau=t-t_g^{\text{up}}+1}^t z_g^\tau \leq x_g^t, \quad \forall g \in G, \quad \forall t \in T \quad (7g)$$

$$\sum_{\tau=t-t_g^{\text{down}}+1}^t z_g^\tau \leq 1 - x_g^t, \quad \forall g \in G, \quad \forall t \in T \quad (7h)$$

$$-R_g^{\text{dn}} \leq p_g^{t+1} - p_g^t \leq R_g^{\text{up}}, \quad \forall g \in G \quad (7i)$$

$$u_g^t, z_g^t, y_g^t \in \{0, 1\}, \quad \forall g \in G, \quad t \in T \quad (7j)$$

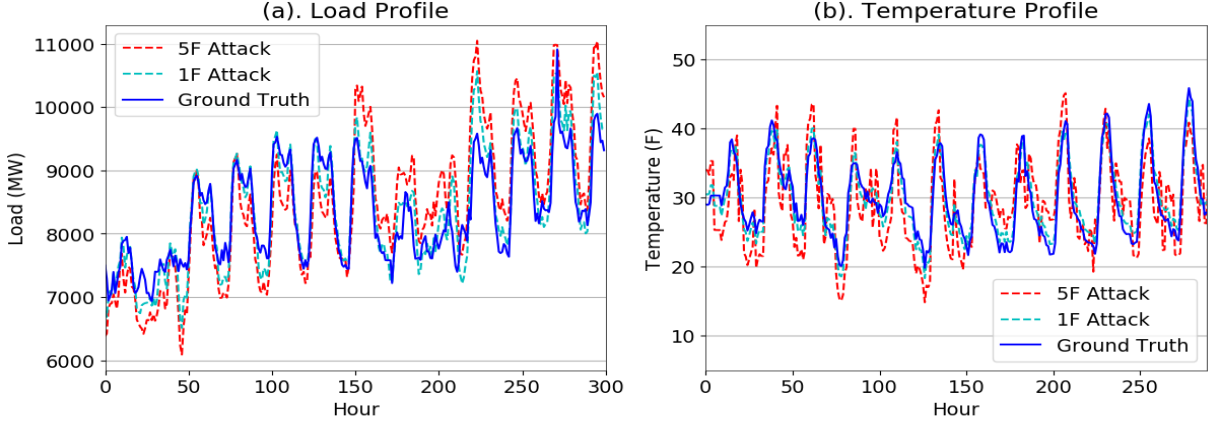


Fig. 2: Visualization of 300 hours forecasts based on original and adversarial temperature data for the aggregated load of Switzerland. The load forecasting model is a recurrent neural networks, while attack perturbations on test dataset are generated by using the gradient estimation method. The attack tries to minimize the load in the first 150 hours and to maximize the load in the latter 150 hours. (a). Load forecasting results; (b). false data injections on temperature.

where $\delta(d)$ is the set of lines connected to node d ; u_g^t is the binary decision variable of the commitment status of generator g at time t , with 1 indicating g is online; p_g^t is the real power output of generator g at time t ; all the u_g^t 's and p_g^t 's are collected together into vectors \mathbf{u} and \mathbf{p} ; $C(\mathbf{p})$ and $S(\mathbf{u})$ represent the dispatch costs and startup and shutdown costs, respectively, of all the generators in all periods; the constraints are system-wide power balance constraint (7b), generation limits constraints (7c), line flow limits (7d), power balance at each node (7e), generator logical constraint (7f), minimum up time constraint (7g), minimum down time constraint (7h) and ramping constraints (7i). We also keep a fixed reserve margin throughout the simulation. Once solved, the operator gets the schedule for the set of online generators G_t at each time t .

The attacker injects $\tilde{\mathbf{X}}$ to a group of compromised load forecasts, such that the nodal load forecasts maliciously change from \hat{L}_d^t to \tilde{L}_d^t . The attacker is not only constrained by the average deviations upon $\tilde{\mathbf{X}}$ (Equation (2)), but also constrained on number of compromised loads N_{adv} that are allowed to inject perturbations throughout the day. From the attacker's perspective, it is then most harmful to find a constrained adversarial day-ahead commitment schedule \tilde{G}_t via \tilde{L}_d^t , which maximizes the solution for the following real-time ED:

$$\min_{\mathbf{p}_t, \mathbf{LS}_t} C(\mathbf{p}_t) + C(\mathbf{LS}_t) \quad (8a)$$

$$\text{s.t. } \sum_{g \in G_t} p_g^t + \sum_{d \in D} LS_d^t = \sum_{d \in D} L_d^t, \quad (8b)$$

$$p_g^{\min} \leq p_g^t \leq p_g^{\max}, \quad \forall g \in \tilde{G}_t \quad (8c)$$

$$f_l^{\min} \leq f_l^t \leq f_l^{\max}, \quad \forall l \in F, \forall t \in T \quad (8d)$$

$$\sum_{g \in \delta(d)} p_g^t + \sum_{k \in \delta(d)} f_k^t = L_d^t, \quad \forall d \in D, t \in T \quad (8e)$$

where from the system operator's perspective, ED aims to find the real power dispatch at time t , \mathbf{p}_t , that minimize the dispatch costs at time t , $C(\mathbf{p}_t)$, considering system-wide power balance constraint (8b), generation limits constraints (8c), line flow limits (8d) and power balance at each node (8e).

B. Attack Strategies

Under normal operating conditions, the load forecasting algorithms provide accurate forecasts on day-ahead load for system operators to solve (7). When the system is under attack, the attacker chooses a group of load buses to inject adversarial temperature forecasts, such that generation schedule coming from day-ahead planning stage is deviating from the normal schedule. Such adversarial generation schedules are likely to cause malicious operation patterns, e.g., increased system costs, load shedding, no feasible generation dispatch or violation of ramping constraints. Essentially, the attacker wants to answer the following questions to find the attacks:

- Which group of load buses should be compromised to inject $\tilde{\mathbf{X}}$?
- How to generate $\tilde{L}_j^t, t = 1, \dots, T$ for compromised load bus j , such that (8) is maximized?

Under the case all system parameters are known, the attacker's optimization problem is tri-level with integer constraints, which is a very challenging problem to solve. Rather than the standard approach through KKT conditions, we design a modified greedy search algorithm for the attacker to find the most vulnerable loads that cause system-level misoperations. As described in Algorithm 1, the attacker follows a modified best-first search algorithm to implement attacks on compromised nodal forecasts iteratively [24]. For each iteration, the attacker checks the lines and generators which are approaching operating limits (e.g., line flow capacity, generation capacity), and finds the most vulnerable load j based on neighboring lines and generators. Then the attacker raises $\tilde{\mathbf{X}}_j$ and $\tilde{\mathbf{X}}_j$ which maximizes and minimizes the load forecasts respectively. By solving (7) using the resulting \tilde{L}_j and \tilde{L}_j , the attacker checks which candidate attack is more prone to make day-ahead commitment schedule \tilde{G}_t different from G_t and keeps it as \tilde{L}_j in this iteration, while j is added to the set of compromised loads $\tilde{\mathbf{L}}$.

When the attacker does not know the parameters of underlying system such as network topology, line flow limits, each generator's capacity and ramp constraints, it is not possible

for the attacker to find the optimal attacks, and the attacker just randomly chooses a set of loads to attack using either $\tilde{\mathbf{X}}_j$ or $\tilde{\mathbf{X}}_j$ for $j \in \tilde{\mathbb{L}}$. In the next section, we will show attacks on real-world load data upon IEEE standard systems that reveal the vulnerabilities brought by either strategic or random load forecasting attacks.

Algorithm 1 Best-First Search

Input: Forecasted loads $L_1, L_2, \dots, L_{|D|} \in \mathbb{R}^{|D| \times T}$

Initialize: Sets of Compromised Loads $\tilde{\mathbb{L}} \leftarrow \emptyset$

Initialize: Maximum Temperature perturbations ϵ

Initialize: Number of compromised loads $k = 0$

Solve UC under normal forecasts

while $G_t \neq \tilde{G}_t, t = 1, \dots, T$ and $k \leq N_{adv}$ **do**

Find the most vulnerable node j and attack direction

Gradient Estimation attacks to find \tilde{L}_j

$\tilde{\mathbb{L}} \leftarrow j$

Solve (7) to get $\tilde{G}_t, k+ = 1$

end while

Solve (8) in real-time with \tilde{G}_t

V. CASE STUDIES

In this section, we show a detailed simulation on real-world Swiss load data, and show the threats posed by our black-box data injection attacks on both the load forecasting algorithm itself and the power system operations. Thorough evaluations on 14-bus system indicate that proposed attacks cause load shedding with very high probability, while tests on 118-bus system indicate by compromising a small portion of nodal load forecasts, the attacker can cause security threats over the whole networks.

A. Experimental Setup

Dataset Description: We collected and queried hourly Swiss actual load data from European Network of Transmission System Operators for Electricity (ENTSO-E)'s API³ ranging from Jan 1st, 2015 to May 16th, 2017. The nominal load values are in the range of $[6,500MW, 9,500MW]$. We followed [25] to collect day-ahead historical weather forecasts coming from major cities in Switzerland such as Zurich, Basel, Lucerne and etc. All the weather data were queried from Dark Sky API⁴. We also collected other indicator features \mathbf{X}^{index} , such as one-hot vectors of hour of day, weekdays and seasons. We evaluated the attack threats on split-out test data, and the forecast model parameters are kept away from the attacker throughout the black-box simulations.

Power Systems Setup: We set up the IEEE 14-bus and IEEE 118-bus to study the system vulnerabilities brought by load forecasting attacks [26]. The grid has a total capacity of 15,500MW and 14,949.3MW respectively, which are both over 1.5 times of yearly peak load. We set the spinning reserve requirement as 3% of the total forecasted demand based on [27]. The models of UC and ED are implemented in

Python using PyPSA [28], and these two modules are directly interfaced with the load forecasting and attack algorithms using Tensorflow [29]. No load shedding occurs when clean load forecasts are used.

Model Training and Attack Implementation: We set up an RNN with 3 layers for day-ahead load forecasting, and use standard stochastic gradient descent methods for model training [30]. Once the validation error converged, our RNN model reports an 1.58% test error in mean absolute percentage error (MAPE), which are comparable to the errors reported in several recent studies on load forecasting [7], [8]. We use L_∞ constraints on the attacker's capability (2b), such that the attacker is constrained by the maximum deviation of perturbed temperature values.

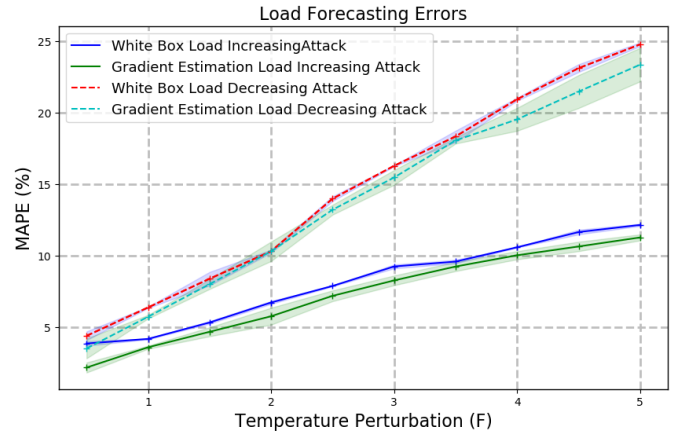


Fig. 3: The forecast MAPE under attacks to increase the load and attacks to decrease the load using white box attacks and gradient estimation attacks. Simulation are run for three times with different random seeds, and shaded area denotes the variance.

B. Load Forecasting Performance

We calibrate and compare the load forecasting model performance with and without adversarial attacks on test datasets. Though the forecasting model exhibits good performances on clean test data, we inject different level of perturbations generated by gradient estimation method, and found the forecasting performance decrease drastically as the adversarial perturbations become larger. In Figure 2 we visualize the RNN's load forecasting results for 300 hours using gradient estimation algorithm with maximum perturbation on temperature of $1F$ and $5F$ respectively. The attacker tries to increase the load in the first 150 hours, and to decrease the load in the latter hours. We observe that the algorithm finds the correct attack direction to either increase or decrease the load. What's more, with only $1F$ deviation on temperatures, the load forecasts changes over 500MW at some time steps. When the attacker increases the perturbation to $5F$, large forecasts error over 1,200MW are observed for the Swiss load. The temperature profile before and after attack still looks similar, which could avoid system operators' security inspection (Figure 2(b)). In Figure 3 we show by either attacking to increase the forecasted loads or attacking to decrease the forecasted loads, the gradient estimation

³<https://transparency.entsoe.eu/>

⁴<https://darksky.net/forecast/47.3769,8.5414/us12/en>

attacks could achieve similar results compared to its white-box counterparts. More essentially, few degrees of malicious perturbations have caused large deviations on forecasted loads.

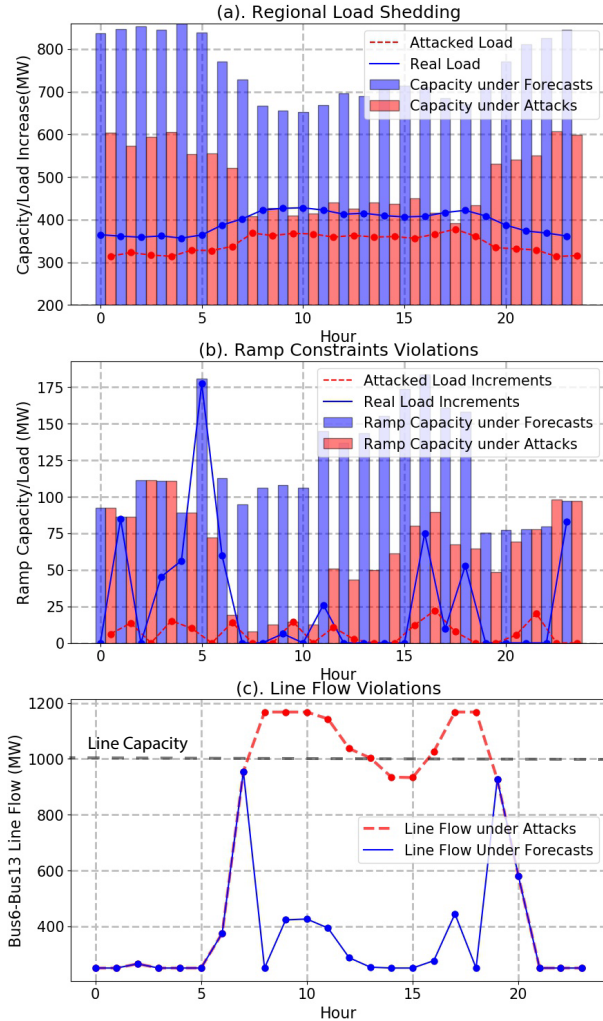


Fig. 4: IEEE 14-bus threats visualization. (a). An example showing that forecasts under attack would cause generation limits violations when real loads exceed total generation capacity; bars indicate generators’ available capacity. (b). An example showing that forecasts under attack would cause violation on ramp constraints during economic dispatch; bars indicate generators’ available total up-ramp capabilities. Maximum allowed perturbations are $4F$. (c). An example showing that line flow exceeds limits.

C. Impacts of Attacks on System Operations

We find small adversarial perturbations over load forecasting input features even cause severe consequences on the power system operations. We are particularly interested in the case when G_t is different from G_t that causes infeasible solution for ED without load shedding. In Figure 4, we visualize different constraint violations when ED is solved.

We ran a thorough evaluation for 122 test days on the 14-bus system, and evaluate if the attackers could cause load shedding under different level of knowledge and capabilities. In Figure 5, we compare the number of days attacker cause load shedding with or without system topology information. In

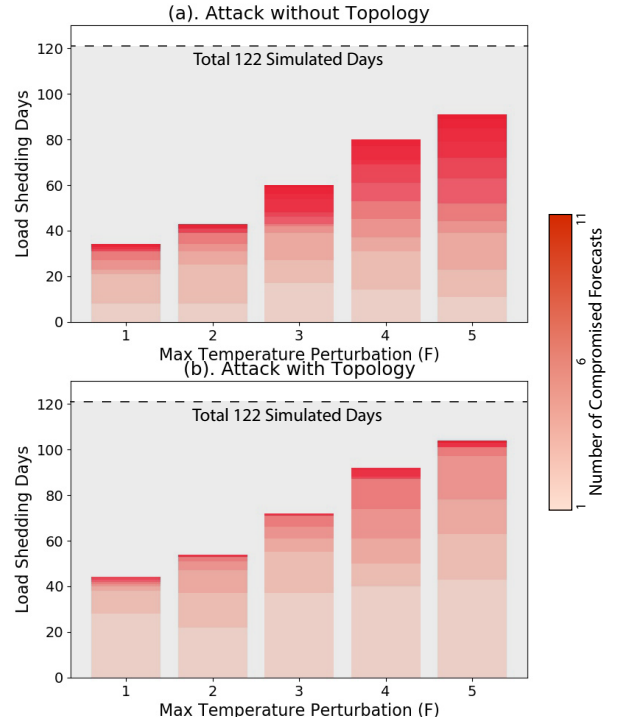


Fig. 5: Load shedding caused by load forecasting attacks on IEEE 14-bus test case. For the 122 test days using Swiss load data, we compare the attack results when attacker either knows or does not know network topology. Larger perturbations on temperatures and extra information on topology helps the attacker to cause load shedding with fewer compromised loads.

both cases, with larger perturbations added to the temperature forecasts, it is more likely to get an infeasible commitment schedule that causes load shedding. At one extreme, when the attacker knows network topology and is able to inject a $5F$ perturbation strategically on compromised load based on Algorithm 1, load shedding occurs on more than 100 of the test days. At the other extreme, the attacker causes over 40 days having shedded loads by only compromising one nodal load forecasts (Figure 5(b)). More surprisingly, when the attacker does not know network topology and just selects compromised load randomly using either \tilde{X}_j or \tilde{X}_j , the system operation is still very vulnerable to the proposed load forecasting attacks (Figure 5(a)).

To further evaluate the attack’s threats brought upon power system operation, one-day attack example on 118-bus test case is illustrated in Figure 6. We assume the attacker is able to design greedy attacks based on known topology information in this case, and show that by only compromising a small subset of nodal load forecasts, the resulting day-ahead commitment schedule under attacks is shutting off several generators compared to that under normal forecasts. The network already observes a series of operation threats such as overflown lines, on-capacity generators and load shedding.

VI. DISCUSSION AND CONCLUSION

In this paper, we studied the potential vulnerabilities generally existing in many load forecasting algorithms. Such vulnerabilities have been overlooked by the development of most if not all forecasting techniques. We designed a data

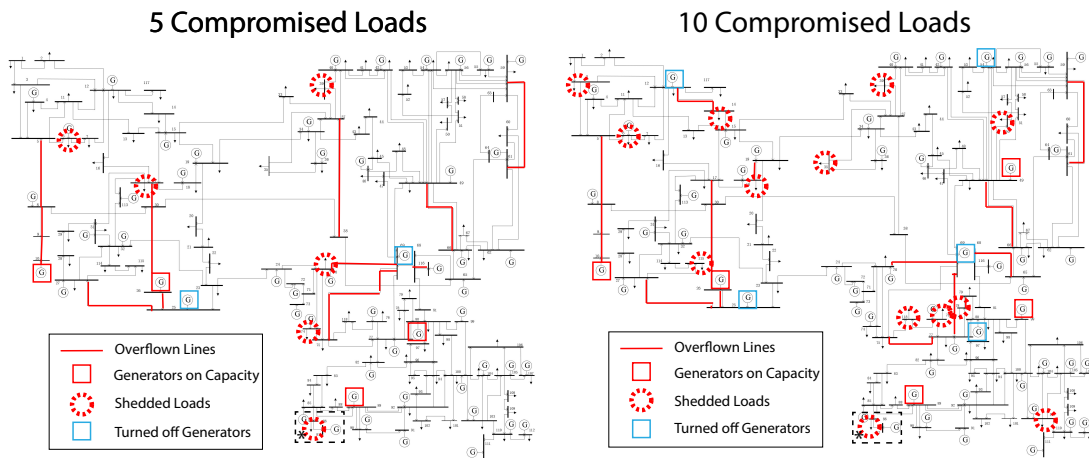


Fig. 6: Threats posed by black-box attack on load forecasting algorithm in the IEEE 118-bus test case. By greedily attacking the nodal day-ahead load forecasts that leads to malicious commitment schedule, the attacker can incur a series of system operations threats such as load shedding, a number of overflow lines and on-capacity generators. For the bottom star-marked load shedded bus, the generator is off in both sub figures.

injection attack which does not require parameters of the forecast algorithms, but leads to large increase in forecast errors. The proposed attack could adversarially impact the decision making process for system operators. Experiments on real-world load datasets demonstrate such threats over power system operations. Such threats model along with damage analysis indicate that there needs more security evaluations in the design and implementation of load forecasting algorithms. In order to mitigate the damages brought by such false data injection attacks, countermeasures such as anomaly detection as well as other robust statistics are strongly recommended.

REFERENCES

- [1] G. Gross and F. D. Galiana, "Short-term load forecasting," *Proceedings of the IEEE*, vol. 75, no. 12, pp. 1558–1573, 1987.
- [2] B. F. Hobbs, S. Jitrapakulsum, S. Konda, V. Chankong, K. A. Loparo, and D. J. Maratukulam, "Analysis of the value for unit commitment of improved load forecasts," *IEEE Transactions on Power Systems*, vol. 14, no. 4, pp. 1342–1348, 1999.
- [3] J. G. De Gooijer and R. J. Hyndman, "25 years of time series forecasting," *International journal of forecasting*, vol. 22, no. 3, pp. 443–473, 2006.
- [4] E. Ceperic, V. Ceperic, A. Baric *et al.*, "A strategy for short-term load forecasting by support vector regression machines," *IEEE Transactions on Power Systems*, vol. 28, no. 4, pp. 4356–4364, 2013.
- [5] J. Contreras, R. Espinola, F. J. Nogales, and A. J. Conejo, "Arma models to predict next-day electricity prices," *IEEE transactions on power systems*, vol. 18, no. 3, pp. 1014–1020, 2003.
- [6] H. S. Hippert, C. E. Pedreira, and R. C. Souza, "Neural networks for short-term load forecasting: A review and evaluation," *IEEE Transactions on power systems*, vol. 16, no. 1, pp. 44–55, 2001.
- [7] W. Kong, Z. Y. Dong, Y. Jia, D. J. Hill, Y. Xu, and Y. Zhang, "Short-term residential load forecasting based on lstm recurrent neural network," *IEEE Transactions on Smart Grid*, 2017.
- [8] F. L. Quilumba, W.-J. Lee, H. Huang, D. Y. Wang, and R. L. Szabados, "Using smart meter data to improve the accuracy of intraday load forecasting considering customer behavior similarities," *IEEE Trans. Smart Grid*, vol. 6, no. 2, pp. 911–918, 2015.
- [9] F. Wang, B. Liu, and T. Hong, "Electric load forecasting with recency effect: A big data approach," *International Journal of Forecasting*, vol. 32, no. 3, pp. 585–597, 2016.
- [10] P. McDaniel and S. McLaughlin, "Security and privacy challenges in the smart grid," *IEEE Security & Privacy*, no. 3, pp. 75–77, 2009.
- [11] Y. Mo and B. Sinopoli, "Secure control against replay attacks," in *Communication, Control, and Computing, 2009. Allerton 2009. 47th Annual Allerton Conference on*. IEEE, 2009, pp. 911–918.
- [12] O. Kosut, L. Jia, R. J. Thomas, and L. Tong, "Malicious data attacks on smart grid state estimation: Attack strategies and countermeasures," in *Smart Grid Communications (SmartGridComm), 2010 First IEEE International Conference on*. IEEE, 2010, pp. 220–225.
- [13] Y. Liu, P. Ning, and M. K. Reiter, "False data injection attacks against state estimation in electric power grids," *ACM Transactions on Information and System Security (TISSEC)*, vol. 14, no. 1, p. 13, 2011.
- [14] G. N. Ericsson, "Cyber security and power system communication essential parts of a smart grid infrastructure," *IEEE Transactions on Power Delivery*, vol. 25, no. 3, pp. 1501–1507, 2010.
- [15] L. Xie, Y. Mo, and B. Sinopoli, "False data injection attacks in electricity markets," in *Smart Grid Communications (SmartGridComm), 2010 First IEEE International Conference on*. IEEE, 2010, pp. 226–231.
- [16] S. Tan, W.-Z. Song, M. Stewart, J. Yang, and L. Tong, "Online data integrity attacks against real-time electrical market in smart grid," *IEEE Transactions on Smart Grid*, vol. 9, no. 1, pp. 313–322, 2018.
- [17] J. Luo, T. Hong, and S.-C. Fang, "Benchmarking robustness of load forecasting models under data integrity attacks," *International Journal of Forecasting*, vol. 34, no. 1, pp. 89–104, 2018.
- [18] Y. Chen, Y. Tan, and B. Zhang, "Exploiting vulnerabilities of load forecasting through adversarial attacks," in *2019 ACM E-Energy Conference*. ACM, 2019.
- [19] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *arXiv preprint arXiv:1312.6199*, 2013.
- [20] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, "The limitations of deep learning in adversarial settings," in *Security and Privacy (EuroS&P), 2016 IEEE European Symposium on*. IEEE, 2016, pp. 372–387.
- [21] Y. Chen, Y. Tan, L. Zhang, and B. Zhang, "Vulnerabilities of power system operations to load forecasting data injection attacks," *arXiv preprint arXiv:1906.04926*, 2019.
- [22] D. S. Kirschen and G. Strbac, *Fundamentals of power system economics*. John Wiley & Sons, 2018.
- [23] A. N. Bhagoji, W. He, B. Li, and D. Song, "Practical black-box attacks on deep neural networks using efficient query mechanisms," in *European Conference on Computer Vision*. Springer, Cham, 2018, pp. 158–174.
- [24] R. Dechter and J. Pearl, "Generalized best-first search strategies and the optimality of a," *Journal of the ACM (JACM)*, vol. 32, no. 3, pp. 505–536, 1985.
- [25] D. L. Marino, K. Amarasinghe, and M. Manic, "Building energy load forecasting using deep neural networks," in *Industrial Electronics Society, IECON 2016-42nd Annual Conference of the IEEE*. IEEE, 2016, pp. 7046–7051.
- [26] R. Christie, "Power systems test case archive," *Electrical Engineering dept., University of Washington*, 2000.

- [27] Y. Rebours and D. Kirschen, "What is spinning reserve," *The University of Manchester*, vol. 174, p. 175, 2005.
- [28] T. Brown, J. Hörsch, and D. Schlachtberger, "PyPSA: Python for Power System Analysis," *Journal of Open Research Software*, vol. 6, no. 4, 2018. [Online]. Available: <https://doi.org/10.5334/jors.188>
- [29] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard *et al.*, "Tensorflow: a system for large-scale machine learning," in *OSDI*, vol. 16, 2016, pp. 265–283.
- [30] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proceedings of COMPSTAT'2010*. Springer, 2010, pp. 177–186.
- [31] N. Papernot, P. McDaniel, and I. Goodfellow, "Transferability in machine learning: from phenomena to black-box attacks using adversarial samples," *arXiv preprint arXiv:1605.07277*, 2016.
- [32] H. Hosseini, Y. Chen, S. Kannan, B. Zhang, and R. Poovendran, "Blocking transferability of adversarial examples in black-box learning systems," *arXiv preprint arXiv:1703.04318*, 2017.
- [33] Y. Chen, Y. Tan, and D. Deka, "Is machine learning in power systems vulnerable?" in *2018 IEEE International Conference on Communications, Control, and Computing Technologies for Smart Grids (Smart-GridComm)*. IEEE, 2018, pp. 1–6.
- [34] K. Chen, K. Chen, Q. Wang, Z. He, J. Hu, and J. He, "Short-term load forecasting with deep residual networks," *IEEE Transactions on Smart Grid*, 2018.
- [35] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Proceedings of the 27th international conference on machine learning (ICML-10)*, 2010, pp. 807–814.

APPENDIX

A. Details on Learn and Attack Algorithm

In addition to the gradient estimation attack introduced in the main text, we also consider the learn and attack setting, where we assume the attacker does not have access to the model parameters, and there is no query access to the model. The only knowledge the attacker has is a historical dataset $\tilde{\mathcal{D}}_{tr}$, which includes same features under same distributions as data set \mathcal{D}_{tr} used to train the load forecasting model⁵. The proposed attack algorithm consists of a *training phase* and an *attack phase* as shown in Figure 7. In the training phase, the attacker trains substitute model $f_{\tilde{\theta}}$ based on $\tilde{\mathcal{D}}_{tr}$ to minimize the training loss. In the attack phase, the attacker pretends that the substitute model is the true load forecast model and performs white-box attacks on it to find the attack vectors. This strategy is based on the assumption that the substitute model behaves similarly to the true model not only for the training data \mathbf{X} , but also for the attack vector $\tilde{\mathbf{X}}$. Then by injecting $\tilde{\mathbf{X}}$ into the true load forecasting model, the forecast values go to attacker's desired directions.

It is useful to evaluate the *transferability* of proposed attacks across different set of models f_{θ} and $f_{\tilde{\theta}}$. The phenomenon of transferability in adversarial attacks for machine learning models have been discussed in [31], [32], where adversarial instance generated using $f_{\tilde{\theta}}$ can be also treated as an adversarial instance by f_{θ} with high probability. The theoretical understanding of why attacks transfer remains an open question and is out of scope for this paper. In Figure 8 we show such *transferability* also exists in the load forecasting model using same test case on Switzerland load forecasting. The temperature inputs are generated by implementing the iterative gradient update based on a substitute model under L_{∞} -norm

⁵In Learn and Attack setting, we make assumption that the attacker know the family of targeted load forecasting model, e.g., a feedforward neural networks or a Recurrent Neural Networks.

of attack perturbations, yet such adversarial temperature values also mislead the (unknown) true load forecasting model to be wildly inaccurate.

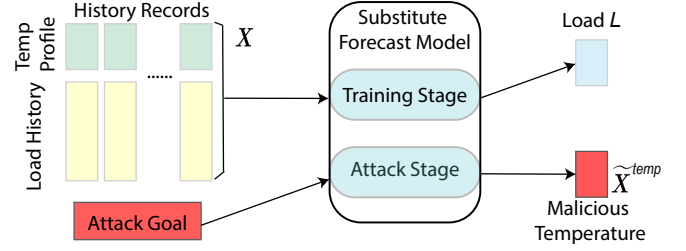


Fig. 7: The attack schematic for learn and attack approach. During the training stage, the attacker uses historical data to learn a substitute forecast model; during attack stage, the attacker finds attack vector using substitute model and transfers it to unknown targeted system.

B. Details on Best First Search Attack Selection

As described in Algorithm 1 in Section IV-B, since it is computationally inefficient and challenging for the attacker to solve a tri-level attack problem (attacker solves an adversarial UC (7) using \tilde{L}_q^t , operator minimizes ED costs (8) using \tilde{G}_t , attacker maximizes ED costs (8) using \tilde{G}_t), we propose a modified best-first search algorithm [24] for attackers to compromise a limited number of nodal forecasts. Essentially, in order to solve the tri-level problem expressively to find attack vectors, our proposed attack on networks iteratively finds the most vulnerable node to inject attacks.

Suppose for $|Q|$ loads in the network, the attacker can at most N_{adv} nodal load forecasts to avoid detection by system operators. Meanwhile, the constraints on attacker's capabilities described in Section III shall hold throughout data perturbation attacks. If an attacker could cause the day-ahead UC schedule $\tilde{G}_t, t = 1, \dots, T$ using $\tilde{L}_q^t, q = 1, \dots, |Q|, t = 1, \dots, T$ shift from $G_t, t = 1, \dots, T$, then it is expected the solution of ED will change. There are several possible circumstances by using \tilde{L}_q^t :

- Increasing the load maliciously will possibly incur extra system costs, such as starting to operate redundant generators, using more expensive generation combinations and etc;
- Decreasing the load maliciously will possibly incur infeasible generation schedules during real-time dispatch, since there may be fewer generators scheduled than normal conditions, which will cause generators reach capacity or line flow exceed limits;
- Decreasing the peak value of load maliciously will possibly cause UC solver ignore peak values, which may cause generators reach ramping limits.

In our proposed algorithm for attackers to find most vulnerable nodal forecasts and inject attack forecasts, we design an iterative search scheme. At each iteration, the attacker solves day-ahead UC (7) based on current \tilde{L} , and check which generator's schedule is most prone to change. Then the attacker decides the next compromised load node j . Since it is more possible to change UC schedule by changing the load

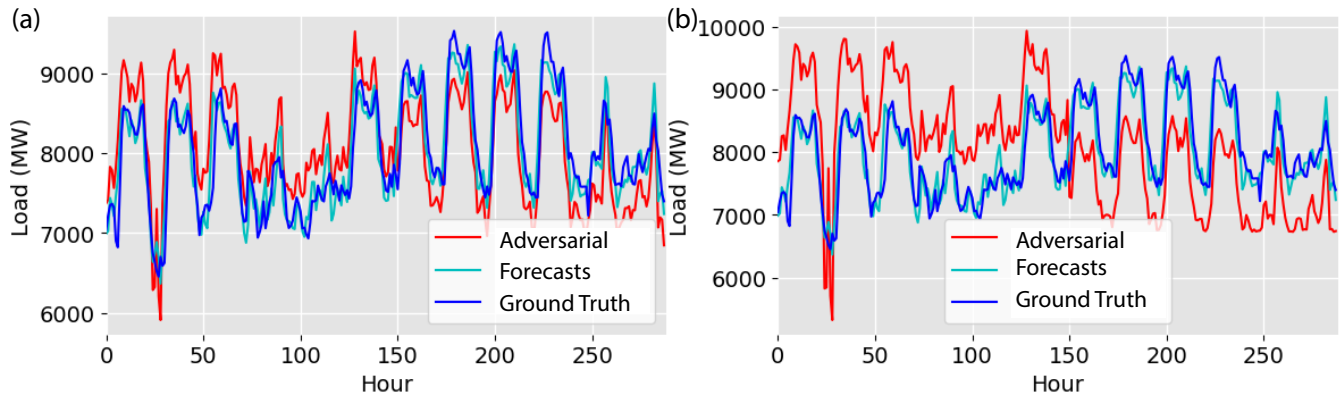


Fig. 8: We show 300 hours forecasts based on original and adversarial temperature data for the aggregated load of Switzerland. The load forecasting algorithm is an recurrent neural networks with inputs composed of past load, regional temperature forecast values and weather indicators. The attack perturbations are generated by using the `learn` and `attack` method, and it implements load maximization strategy in the first 150 hours and load minimization strategy in the latter 150 hours. (a). Load forecasting results with temperature attack constraint of (maximum perturbations) $1F$; (b). load forecasting results with temperature attack constraint of $5F$.

profile to greater extents, the attacker then proposes two attack samples $\tilde{\mathbf{X}}_j$ and $\tilde{\mathbf{X}}_j$ which minimizes and maximizes the nodal forecasts respectively. Such adversarial nodal load forecast is inserted into next iteration’s UC problem.

In our implementation, we simply find the prone-to-change generation schedule based on following criteria:

- The generator’s dispatch during day-ahead UC is either approaching to generation limit or reaching zero output;
- The line flow during day-ahead UC is approaching line flow limits;

And simulation results indicate such search algorithm is quite efficient to find the attack injections. Other criteria such as ramping constraints and generation on-off limits could be also utilized to find the most vulnerable nodal forecasts.

For the case when system topology and parameters are unknown, the attacker just selects a random set of compromised load and inject attacks perturbations using either $\tilde{\mathbf{X}}_j$ or $\tilde{\mathbf{X}}_j$ for the compromised nodal loads. The result shown in Fig. 5 validates that in the case without topology knowledge, compromising the same number of loads are causing fewer load shedding days than in the case when topology is known.

C. Details on Attacks Implementation

In addition to the results we have shown in the main texts using RNN as load forecasting algorithm, we added ablation study on different load forecasting algorithms, as well as attack performance and computation time analysis.

1) Forecasting Method: Feed-Forward Neural Networks

A multi-layered, feed-forward neural networks (NN) has been widely used to represent the nonlinearities between input features and output forecasts [6]. For the input layer of neural networks, each neuron represents one feature of training input, and all features of past H steps ($\mathbf{X}_{t-H}, \dots, \mathbf{X}_t$) are stacked as the inputs. For each intermediate layer, NN could have a tunable number of hidden units, which represent the input feature combinations.

Recurrent Neural Networks As described in main texts, RNN feeds each step’s input \mathbf{X}_t sequentially, and outputs a hidden unit to represent the feature combination of current input and historical features. The last neuron outputs the forecasted load values in the load forecasting attack [33].

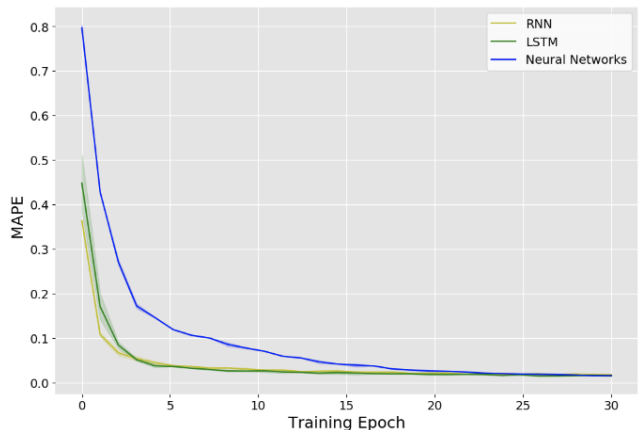


Fig. 9: All three forecasting models, show convergence of forecast error on validation data as training evolves. Shaded areas show the variance of MAPE.

Long Short-Term Memory Long Short-Term Memory network (LSTM) is designed to deal with the vanishing gradient problem existing in the RNN with long-time dependencies [7], [34]. The major improvements over RNN are the design of “forget” gates to model the temporal dependencies and capture long time dependencies in load patterns more accurately.

2) *Training and Attack Details:* We set up all load forecasting models using Tensorflow [29] package in Python. Standard model architectures such as Dropout layers and nonlinear activation functions (e.g., ReLU or Sigmoid functions) are adopted in the deep learning models [35]. Since all three networks are set up to solve the load forecasting regression problem, we set the first layer having most neurons, and decrease the number of

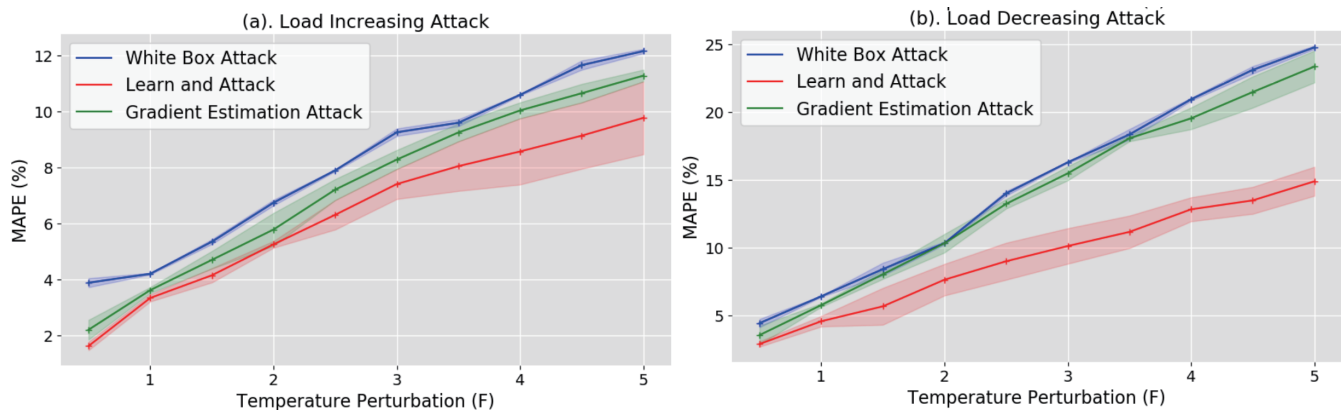


Fig. 10: The forecast MAPE under (a). attacks to increase the load; and (b). attacks to decrease the load. Simulations are run for three times with different random seeds for all three attacks, and shaded area denotes the variance.

Forecasts Models	NN	RNN	LSTM
Number of Layers	4	3	3
Training Epochs	20	30	30
Hidden Units in First Layer	512	64	64

TABLE I: Model architectures and training configurations for load forecasting algorithms used in the simulations.

Forecasts Error (MAPE)	Clean Data	Learn and Attack	Gradient Estimation
NN	1.68%	12.72%	13.09%
RNN	1.58%	9.82%	11.68%
LSTM	1.51%	9.04%	11.87%

TABLE II: Forecasts errors evaluated on clean test data and adversarial data for 3 different forecast models. Allowed maximum perturbations are $4F$.

units in subsequent layers. Table I records our load forecasting model setup. We split 80% of data as our training sets, and use the remaining 20% of data on validating and evaluating the load forecasting prediction accuracy, attack performance and case studies on market operations. Such data collection procedures could also be applied in an online fashion so that attacker could inject real-time adversarial attacks into certain load forecasting models.

As shown in Figure 9, all three load forecasting algorithms' validation loss are converged during training, and we use the trained model in the subsequent planning and operation problem as well as the testbed for attack algorithms. Plots are showing the mean and variance during 3 runs.

For the substitute model training of learn and attack method, we keep the training set \mathcal{D} same as the load forecasting model training set \mathcal{D} . Decreasing the size of \mathcal{D} or using different substitute dataset could decrease the performance of learn and attack. Table II compares all three load forecasting models' performance using clean and adversarial data. For both learn and attack and gradient estimation algorithms, they distort all three load forecasting models' output and increase model's forecast error. Gradient estimation attack works generally better for all three models, and this is due to estimating the gradients via querying f_θ directly is more accurate than calculating it

from the substitute model and transferring to f_θ .

In Figure 10, we evaluate RNN's load forecasting performance under two attack strategies: load maximization or load minimization. We observe gradient estimation attack causes similar MAPE compared to white box attack. The load decreasing attack is normally more successful than load increasing attack in terms of MAPE. Load minimization attack is more harmful results than load increasing ones, since increased forecasts only let system operators start up more generations, while adversarially decreasing the forecasted load leads to wrong generation decisions that fails to meet the larger real load.

Forecasts Models	NN	RNN	LSTM
Training Time	12.988	47.998	143.830
Learn and Attack	0.133	0.157	0.579
Gradient Estimation Attack	0.082	0.119	0.253

TABLE III: Computation time (in seconds) for load forecasting model training and implementation time for attacks.

3) *Computation Time*: We recorded the computation time for neural network training and the implementation time for two proposed attack algorithms. All time are recorded on a laptop with Intel 2.3GHz Core i5-8259U 4 Cores CPU and 8 GB RAM. The training time for NN, RNN and LSTM are calculated for 20, 30 and 30 epochs respectively. The implementation time for the attacks are averaged over all test instances. We observed that learn and attack approach takes longer time than gradient estimation due to the longer time taken to calculate gradient signs over the whole neural networks; and as LSTM includes more complicated model architectures, it takes longer time to find the adversarial instance. Yet compared to the long model training time and application scenarios of day-ahead forecasts, the attacker is still efficient enough to find the adversarial perturbations. Such efficient computation enables attacker to find the most vulnerable loads or to attack the system operations in very short time.

4) *Code Availability*: The implementation code for forecasts, attacks, and power system operations are all available at https://github.com/chennnnyize/load_forecasts_attack.