

Verification of Neural Network Behaviour: Formal Guarantees for Power System Applications

Andreas Venzke, *Student Member, IEEE*, and Spyros Chatzivasileiadis, *Senior Member, IEEE*

Abstract—This paper presents for the first time, to our knowledge, a framework for verifying neural network behavior in power system applications. Up to this moment, neural networks have been applied in power systems as a black box; this has presented a major barrier for their adoption in practice. Developing a rigorous framework based on mixed-integer linear programming, our methods can determine the range of inputs that neural networks classify as safe or unsafe, and are able to systematically identify adversarial examples. Such methods have the potential to build the missing trust of power system operators on neural networks, and unlock a series of new applications in power systems. This paper presents the framework, methods to assess and improve neural network robustness in power systems, and addresses concerns related to scalability and accuracy. We demonstrate our methods on the IEEE 9-bus, 14-bus, and 162-bus systems, treating both N-1 security and small-signal stability.

Index Terms—Neural networks, mixed-integer linear programming, security assessment, small-signal stability.

I. INTRODUCTION

A. Motivation

MACHINE learning, such as decision trees and neural networks, has demonstrated significant potential for highly complex classification tasks including the security assessment of power systems [1]. However, the inability to anticipate the behavior of neural networks, which have been usually treated as a black box, has been posing a major barrier in their application in safety-critical systems, such as power systems. Recent works (e.g. [2]) have shown that neural networks that have high prediction accuracy on unseen test data can be highly vulnerable to so-called adversarial examples (small input perturbations leading to false behaviour), and that their performance is not robust. To our knowledge, the robustness of neural networks has not been systematically evaluated in power system applications before. This is the first work that develops provable formal guarantees of the behavior of neural networks in power system applications. These methods allow to evaluate the robustness and improve the interpretability of neural networks and have the potential to build the missing trust of power system operators in neural networks, enabling their application in safety-critical operations.

B. Literature Review

Machine learning algorithms including neural networks have been applied in power system problems for decades.

A. Venzke and S. Chatzivasileiadis are with the Department of Electrical Engineering, Technical University of Denmark, 2800 Kgs. Lyngby, Denmark e-mail: {andven, spchatz}@elektro.dtu.dk.

This work is supported by the multiDC project, funded by Innovation Fund Denmark, Grant Agreement No. 6154-00020.

For a comprehensive review, the interested reader is referred to [1], [3] and references therein. A recent survey in [4] reviews applications of (deep) reinforcement learning in power systems. In the following, we focus on applications related to power system operation and security. Recent examples include [5] which compares different machine learning techniques for probabilistic reliability assessment and shows significant reductions in computation time compared to conventional approaches. Neural networks obtain the highest predictive accuracy. The works in [6], [7] rely on machine learning techniques to learn local control policies for distributed energy resources in distribution grids. Ref. [8] uses machine learning to predict the result of outage scheduling under uncertainty, while in [9] neural networks rank contingencies for security assessment. Neural networks are used in [10] to approximate the system security boundary and are then included in the security-constrained optimal power flow (OPF) problem.

Recent developments in deep learning have sparked renewed interest for power system applications [11]–[16]. There are a range of applications for which deep learning holds significant promise including online security assessment, system fault diagnosis, and rolling optimization of renewable dispatch as outlined in [11]. Deep neural networks are used in [12] to predict the line currents for different N-1 outages. By encoding the system state inside the hidden layers of the neural network, the method can be applied to unseen N-2 outages with high accuracy [13]. The work in [14] proposes a deep autoencoder to reduce the high-dimensional input space of security assessment and increase classifier accuracy and robustness. Our work in [15] represents power system snapshots as images to take advantage of the advanced deep learning toolboxes for image processing. Using convolutional neural networks, the method assesses both N-1 security and small signal stability at a fraction of the time required by conventional methods. The work in [16] uses convolutional neural networks for N-1 contingency screening of a large number of uncertainty scenarios, and reports computational speed-ups of at least two orders of magnitude. These results highlight the potential of data-driven applications for online security assessment. The black box nature of these tools, however, presents a major obstacle towards their application in practice.

C. Main Contributions

In power system applications, to the best of our knowledge, the robustness of learning approaches based on neural networks has so far not been systematically investigated. Up to now, neural network performance has been evaluated by splitting the available data in a training and a test set,

and assessing accuracy and other statistical metrics on the previously unseen test set data (e.g. [12]–[15]). Recent works (e.g. [2]) have shown that neural networks that have high prediction accuracy on unseen test data can be highly vulnerable to adversarial examples and the resulting prediction accuracy on adversarially crafted inputs is very low [17]. Adversarial examples also exist in power system applications as demonstrated in [18]. This highlights the importance to develop methodologies which allow to systematically evaluate and improve the robustness of neural networks.

In this work, we present for the first time a framework to obtain formal guarantees of neural network behaviour for power system applications, building on recent advances in the machine learning literature [19]. Our contributions are threefold. First, we evaluate the robustness of neural networks in power system applications, through a *systematic* identification of adversarial examples (small input perturbations which lead to false behaviour) or by *proving* that no adversarial examples can exist for a continuous range of neural network inputs. Second, we improve the interpretability of neural networks by obtaining provable guarantees about how continuous input regions map to specific classifications. Third, using systematically identified adversarial examples, we retrain neural networks to improve robustness.

In the rest of this paper, we refer to verification as the process of obtaining formal guarantees for neural network behaviour. These formal guarantees are in the form of continuous input regions in which the classification of the neural network provably does not change, i.e., no adversarial examples exist. Being able to determine the continuous *range* of inputs (instead of discretized samples) that the neural network classifies as safe or unsafe makes the neural network interpretable. Accuracy is no longer a pure statistical metric but can be supported by provable guarantees of neural network behavior. This allows operators to either build trust in the neural network and use it in real-time power system operation, or decide to retrain it. Increasing the robustness, and provably verifying target properties of these algorithms is therefore a prerequisite for their application in practice. To this end, the main contributions of our work are:

- 1) Using mixed-integer linear programming (MILP), we present a neural network verification framework for power system applications which, for continuous ranges of inputs, can guarantee if the neural network will classify them as safe or unsafe.
- 2) We present a systematic procedure to identify adversarial examples and determine neural network input regions in which *provably* no adversarial examples exist.
- 3) We improve the robustness of neural networks by re-training them on enriched training datasets that include adversarial examples identified in a systematic way.
- 4) Formulating the verification problem as a mixed-integer linear program, we investigate, test and apply techniques to maintain scalability; these involve bound tightening and weight sparsification.
- 5) We demonstrate our methods on the IEEE 9-bus, 14-bus, and 162-bus system, treating both N-1 security and small-signal stability. For the IEEE 9-bus system, we

re-train the neural network using identified adversarial examples, and show improvements both in accuracy and robustness.

This work is structured as follows: In Section II, we describe the neural network architecture and training, in Section III we formulate the verification problems as mixed-integer programs and in Section IV we address tractability. In Section V we define our simulation setup and present results on formal guarantees for a range of power system security classifiers.

II. NEURAL NETWORK ARCHITECTURE AND TRAINING

A. Neural Network Structure

Before moving on with the formulation of the verification problem, in this section we explain the general structure of the neural networks we consider in this work [20]. An illustrative example is shown in Fig. 1. A neural network is defined by the number K of its fully-connected hidden layers, with each layer having N_k number of neurons (also called nodes or hidden units), with $k = 1, \dots, K$. The input vector is denoted with $\mathbf{x} \in \mathbb{R}^{N_0}$, and the output vector with $\mathbf{y} \in \mathbb{R}^{N_{K+1}}$. The input to each layer $\hat{\mathbf{z}}_{k+1}$ is a linear combination of the output of the previous layer, i.e. $\hat{\mathbf{z}}_{k+1} = \mathbf{W}_k \mathbf{z}_k + \mathbf{b}_k$, where \mathbf{W}_k is a $N_{k+1} \times N_k$ weight matrix and \mathbf{b}_k is a $N_{k+1} \times 1$ bias vector between layers k and $k+1$. Each neuron in the hidden layers incorporates an activation function, $z = f(\hat{z})$, which usually applies a non-linear transformation to the scalar input. There is a range of possible activation functions, such as the sigmoid function, the hyperbolic tangent, the Rectifier Linear Unit (ReLU), and others. Recent advances in computational power and machine learning have made possible the successful training of deep neural networks [21]. The vast majority of such networks use ReLU as the activation function as this has been shown to accelerate their training [22]. For the rest of this paper we will focus on ReLU as the chosen activation function. A framework for neural network verification considering general activation functions such as the sigmoid function and the hyperbolic tangent is proposed in [23]. ReLU is a piecewise linear function defined as $z = \max(\hat{z}, 0)$. For each of the hidden layers we have:

$$\mathbf{z}_k = \max(\hat{\mathbf{z}}_k, 0) \quad \forall k = 1, \dots, K \quad (1)$$

$$\hat{\mathbf{z}}_{k+1} = \mathbf{W}_{k+1} \mathbf{z}_k + \mathbf{b}_{k+1} \quad \forall k = 0, 1, \dots, K-1 \quad (2)$$

where $\mathbf{z}_0 = \mathbf{x}$, i.e. the input vector. Throughout this work, the max operator on a vector $\hat{\mathbf{z}}_k \in \mathbb{R}^{N_k}$ is defined element-wise as $\mathbf{z}_k^n = \max(\hat{\mathbf{z}}_k^n, 0)$ with $n = 1, \dots, N_k$. The output vector is then obtained as follows:

$$\mathbf{y} = \mathbf{W}_{K+1} \mathbf{z}_K + \mathbf{b}_{K+1} \quad (3)$$

In this work, we will focus on classification networks, that is each of the output states y_i corresponds to a different class. For example, within the power systems context, each operating point \mathbf{x} can be classified as $y_1 = \text{safe}$ or $y_2 = \text{unsafe}$ (binary classification). The input vector \mathbf{x} encodes the necessary information to determine the operating point, e.g. the generation dispatch and loading in the DC optimal power flow (OPF). Each input is classified to the category that corresponds to the largest element of the output vector \mathbf{y} . For example, if $y_1 > y_2$ then input \mathbf{x} is safe, otherwise unsafe.

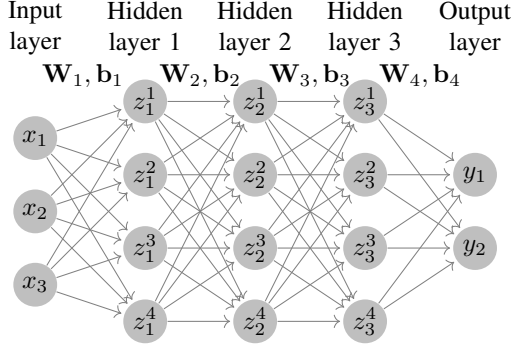


Fig. 1. Illustrative neural network for binary classification: Fully connected neural network with three inputs \mathbf{x} and three hidden layers. Between each layer, a weight matrix \mathbf{W} and bias \mathbf{b} is applied. Each hidden layer has four neurons with non-linear ReLU activation functions. Based on the comparison of the two outputs \mathbf{y} , the input is classified in one of the two categories, i.e. $y_1 > y_2$ or $y_1 < y_2$.

B. Neural Network Training

The training of neural networks requires a dataset of labeled samples \mathbf{S} , where each sample consists of the input \mathbf{x} and the true output classification $\bar{\mathbf{y}}$. The dataset \mathbf{S} is split into a training and a testing dataset denoted with $\mathbf{S}_{\text{train}}$ and \mathbf{S}_{test} , respectively. During training, the weights \mathbf{W} and biases \mathbf{b} are optimized using the training dataset $\mathbf{S}_{\text{train}}$ with respect to a defined objective function. Different objective functions (also called loss functions) exist [20]. In this work we use one of the most commonly employed for classification: the softmax cross entropy, which is defined as follows. First, the softmax function takes the vector of the neural network output \mathbf{y} and transforms it to an equal-size vector of real numbers in the range between 0 and 1, the sum of which have to be equal to 1. This corresponds to a probability distribution of the output of the neural network belonging to a certain class. The probability p_i of the input belonging to the different classes $i \in N_{K+1}$ is defined as:

$$p_i = \frac{e^{y_i}}{\sum_{j=1}^{N_{K+1}} e^{y_j}}, \quad \forall i \in N_{K+1} \quad (4)$$

Note that $e^{\{\cdot\}}$ refers to the exponential function here. We define the softmax cross entropy function as:

$$f(\bar{\mathbf{y}}, \mathbf{p}) = - \sum_{i=1}^{N_{K+1}} \bar{y}_i \log(p_i) \quad (5)$$

This objective function can be understood as the squared error of the distance from the true classification $\bar{\mathbf{y}}$ to the probability distribution over the classes \mathbf{p} predicted by the neural network. This formulation has theoretical and practical advantages over penalizing the squared error directly [24]. During training, we solve the following optimization problem:

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{b}, \mathbf{x}, \mathbf{y}, \bar{\mathbf{y}}, \mathbf{z}, \hat{\mathbf{z}}, \mathbf{p}} \quad & f(\bar{\mathbf{y}}, \mathbf{p}) \\ \text{s.t.} \quad & (1), (2), (3), (4) \text{ and } (5) \end{aligned} \quad (6)$$

The training follows these steps: First, we initialize the weights \mathbf{W} and biases \mathbf{b} randomly. Second, we substitute the variables \mathbf{x} and $\bar{\mathbf{y}}$ with the samples from the training dataset $\mathbf{S}_{\text{train}}$. Based on the current value of weights \mathbf{W} and biases \mathbf{b} , we compute the corresponding neural network prediction \mathbf{y} , and the loss

function (5). Third, using back-propagation, we compute the gradients of the loss function with respect to weights \mathbf{W} and biases \mathbf{b} and update these using stochastic gradient descent (in our simulation studies we use the Adam optimizer [25]). Then, we return to the second step, and repeat this procedure until a defined number of training iterations, also called epochs, is reached. After the training has terminated, we evaluate the neural network performance by calculating its accuracy on the unseen test dataset \mathbf{S}_{test} . This gives us a measure of the generalization capability of the neural network. For a detailed explanation of neural network training please refer to e.g. [26].

III. VERIFICATION AS A MIXED-INTEGER PROGRAM

As the accuracy on the test dataset is not sufficient to provide provable guarantees of neural network behavior, in the following, we formulate verification problems that allow us to verify target properties of neural networks and identify adversarial examples in a rigorous manner. To this end, we first reformulate trained neural networks as mixed-integer programs. Then, we introduce verification problems to a) prove the absence of adversarial examples, b) evaluate the neural network robustness, and c) compute largest input regions with same classification. Finally, we discuss how to include additional constraints on the input to the neural network and how to extend the framework to regression problems.

A. Reformulation of ReLU as Mixed-Integer Program (MIP)

First, we reformulate the ReLU function $\mathbf{z}_k = \max(\hat{\mathbf{z}}_k, 0)$, shown in (1), using binary variables $\mathbf{r}_k \in \{0, 1\}^{N_k}$, following the work in [19]. For all $k = 1, \dots, K$, it holds:

$$\mathbf{z}_k = \max(\hat{\mathbf{z}}_k, 0) \Rightarrow \begin{cases} \mathbf{z}_k \leq \hat{\mathbf{z}}_k - \hat{\mathbf{z}}_k^{\min} (1 - \mathbf{r}_k) & (8a) \\ \mathbf{z}_k \geq \hat{\mathbf{z}}_k & (8b) \\ \mathbf{z}_k \leq \hat{\mathbf{z}}_k^{\max} \mathbf{r}_k & (8c) \\ \mathbf{z}_k \geq \mathbf{0} & (8d) \\ \mathbf{r}_k \in \{0, 1\}^{N_k} & (8e) \end{cases}$$

The entries of the binary vector \mathbf{r}_k signal whether the corresponding ReLU unit in layer k is active or not. Note that the operations in (8a)–(8d) are performed element-wise. For a ReLU unit n , if $r_k^n = 0$, the ReLU is inactive and z_k^n is constrained to be zero through (8c) and (8d) given that the expression $0 \leq \hat{z}_k^n - \hat{z}_k^{\min}$ is true. Conversely, if r_k^n is 1, then z_k^n is constrained to be equal to \hat{z}_k^n through (8a) and (8b) given that the expression $\hat{z}_k^n \leq \hat{z}_k^{\max}$ is true. To ensure that both expressions are always true, the bounds on the ReLU output \hat{z}^{\min} and \hat{z}^{\max} have to be chosen sufficiently large to not be binding, but as low as possible to provide tight bounds in the branch-and-bound algorithm during the solution of the MIP. In Section IV-1, we present a possible approach to tighten these bounds using interval arithmetic (IA).

B. Formulating Verification Problems

Using (8a)–(8e), we can now verify neural network properties by formulating mixed-integer linear programs (MILPs). Without loss of generality, in this paper we assume that the

neural network classifies the output in only two categories: y_1 and y_2 (e.g. safe and unsafe). Note that in case the neural network outputs are of the same magnitude ($y_1 = y_2$), we classify the input as ‘unsafe’; this avoids the risk of classifying a true ‘unsafe’ operating point as ‘safe’ in those cases.

a) Verifying against adversarial examples: Assume a given input \mathbf{x}_{ref} is classified as y_1 , i.e. for the neural network output holds $y_1 > y_2$. Adversarial examples are instances close to the original input \mathbf{x}_{ref} that result to a different (wrong) classification [2]. Imagine for example an image of 2000×1000 pixels showing a green traffic light. An adversarial example exists if by changing just 1 pixel at one corner of the image, the neural network would classify it as a red light instead of the initial classification as green (e.g. see Fig. 16 of [27]). Machine learning literature reports on a wide range of such examples and methods for identifying them, as they can be detrimental for safety-critical application (such as autonomous vehicles). Due to the nature of this problem, however, most of these methods rely on heuristics. Verification can be a very helpful tool towards this objective as it can help us discard areas around given inputs by providing guarantees that no adversarial example exists [19].

Turning back to our problem, assume that the system operator knows that within distance ϵ from the operating point \mathbf{x}_{ref} the system remains safe ($y_1 > y_2$). If we can guarantee that the neural network will output indeed a safe classification for any input $\|\mathbf{x} - \mathbf{x}_{\text{ref}}\| \leq \epsilon$, then we can provide the operator with the guarantees required in order to deploy methods based on neural networks for real-time power system applications. To this end, we can solve the following mixed-integer program:

$$\min_{\mathbf{x}, \hat{\mathbf{z}}, \mathbf{z}, \mathbf{y}} y_1 - y_2 \quad (9a)$$

$$\text{s.t.} \quad (2), (3), (8a) - (8e) \quad (9b)$$

$$\|\mathbf{x} - \mathbf{x}_{\text{ref}}\|_* \leq \epsilon \quad (9c)$$

If the resulting objective function value is strictly positive (and assuming zero optimality gap), then $y_1 > y_2$ for all $\|\mathbf{x} - \mathbf{x}_{\text{ref}}\|_* \leq \epsilon$, and we can guarantee that no adversarial input exists within distance ϵ from the given point. The norm in (9c) can be chosen to be e.g. ∞ -norm, 1-norm or 2-norm. In the following, we focus on the ∞ -norm, which allows us to formulate the optimization problem (9) as MILP. In addition, the ∞ -norm has a natural interpretation of allowing each input dimension in magnitude to vary at most by ϵ . By considering both the 1- and ∞ -norm in (9c), adversarial robustness with respect to all l_p -norms with $p \geq 2$ can be achieved [28]. Conversely, in case input \mathbf{x}_{ref} was originally classified as y_2 , we minimize $y_2 - y_1$ in objective function (9a) instead. Note that we model \mathbf{z} , $\hat{\mathbf{z}}$ and \mathbf{y} as optimization variables in the mixed-integer program in (9), as we are searching for an input \mathbf{x} with output \mathbf{y} that changes the classification. For fixed input \mathbf{x} these optimization variables are then uniquely determined.

b) Adversarial robustness: Instead of solving (9) for a single ϵ and a single operating point \mathbf{x}_{ref} , we can solve a series of optimization problems (9) for different values of ϵ and different operating points \mathbf{x}_{ref} and assess the adversarial accuracy as a measure of neural network robustness. The adversarial accuracy is defined as the share of samples that

do not change classification from the correct ground-truth classification within distance ϵ from the given input \mathbf{x}_{ref} . In our simulation studies, we use all samples in the training data or unseen test data set to evaluate the adversarial accuracy. As such, the adversarial accuracy for a distance measure of zero ($\epsilon = 0$) is equal to the share of correctly predicted samples, i.e. the prediction accuracy. The adversarial accuracy can be used as an indicator for the robustness/brittleness of the neural network: low adversarial accuracy for very small ϵ is in most cases an indicator of poor neural network performance [29].

Furthermore, utilizing this methodology, we can systematically identify adversarial examples to evaluate the adversarial robustness and re-train neural networks to improve their adversarial robustness. First, for a range of different values of ϵ and using either training or test dataset, we solve the optimization problems in (9). For all samples \mathbf{x}_{ref} for which an input perturbation exists within a distance of ϵ that does change the classification, we need to assess whether this change occurs as the sample is located at the power system security boundary (and the input perturbation places the adversarial input across the boundary) or if it is in fact an adversarial example and the change in classification indicates an incorrect boundary prediction by the neural network. This can be achieved by computing the ground-truth classification $\bar{\mathbf{y}}$ for the perturbed input, e.g., by evaluating the system security and stability using conventional methods for the identified potential adversarial inputs, and comparing it to the neural network prediction. The share of identified adversarial examples (i.e. false classification changes) serves as measure of adversarial robustness.

In our simulation studies in Section V, we compute the adversarial accuracy and identify adversarial examples for several illustrative test cases using the proposed methodology. In case the neural network robustness is not satisfactory, we can use the aforementioned procedure to systematically identify adversarial examples from the training dataset, add these to the training dataset and re-train the neural network by optimizing (6)–(7) using the enriched training dataset to improve robustness. Other directions for improving the performance of the neural network include to (i) use a dataset creation method that provides a more detailed description of the security boundary e.g. using the algorithm in [30], and (ii) re-train the neural networks to be adversarially robust by modifying the objective function (5) as outlined e.g. in [29].

c) Computing largest regions with same classification:

To establish trust of neural networks among power system operators it is crucial to be able to determine the range of inputs a neural network will classify as safe or unsafe. The formulation (10) does that by computing the maximum input range around a given input \mathbf{x}_{ref} for which the classification will not change. Note that we achieve this by computing the *minimum* distance to a sample which *does* change the classification:

$$\min_{\mathbf{x}, \hat{\mathbf{z}}, \mathbf{z}, \mathbf{y}, \epsilon} \epsilon \quad (10a)$$

$$\text{s.t.} \quad (2), (3), (8a) - (8e) \quad (10b)$$

$$\|\mathbf{x} - \mathbf{x}_{\text{ref}}\|_* \leq \epsilon \quad (10c)$$

$$y_2 \geq y_1 \quad (10d)$$

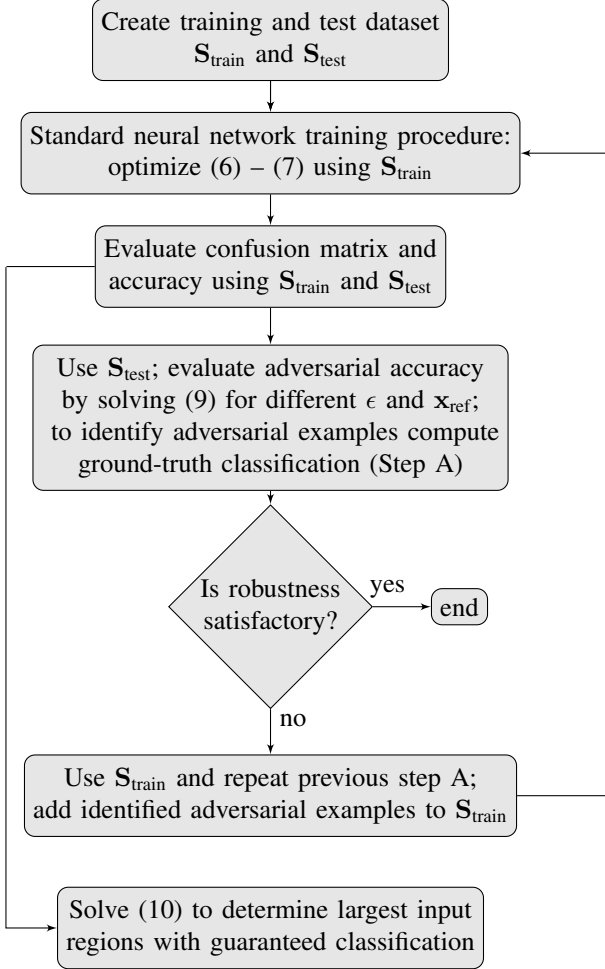


Fig. 2. Flowchart illustrating the methodology: First, we create datasets which are split into training and test set. Using the training set only, we train the neural network. Following standard procedure, we evaluate the neural network performance with the confusion matrix. Then, using the mixed-integer linear reformulation of the trained neural network and the test dataset, we evaluate the adversarial accuracy and identify adversarial examples. If the neural network robustness is not satisfactory, we repeat this step with the training test set, add the identified adversarial examples to the training set and re-train the neural network. Note that we cannot use the test set for this step, as the test set information should not be used in the training process; instead the test set should only be used to evaluate the generalization capability on unseen data. Finally, as shown in the last block as a separate stream, with our methods we are also able to determine input regions around selected input samples with guaranteed classification to provide formal guarantees for neural network behaviour.

Here again we select the ∞ -norm in (10c), turning (10) to a MILP. If input \mathbf{x}_{ref} is classified as y_2 , then we replace (10d) with $y_1 \geq y_2$ instead. Solving verification problem (10) enables us to provide the operator with guarantees that e.g. all system states where Generator 1 produces between 50 MW and 100 MW and Generator 2 produces between 30 MW and 80 MW, will be classified by the neural network as safe. Thus, the neural network is no longer a black box, and it is up to the operator to determine if its results are accurate enough. If they are not, the operator can retrain the neural network with adjusted parameters to increase the accuracy for regions the operator is mostly interested. Alternatively, the operator can follow a risk-based approach attaching a certain confidence/risk level to the neural network results, or the operator can choose to use the neural network for a

limited range of inputs where its results are provably 100% accurate. Our overall methodology is illustrated in a flowchart in Figure 2.

d) *Input constraints*: In all verification problems, we can add additional constraints characterizing the admissible region for the input vector \mathbf{x} and, thus, can obtain larger regions with the same classification. Given that in neural network training it is standard practice to normalize the input between 0 and 1, we add constraint (11) to both formulations in (9) and (10):

$$0 \leq \mathbf{x} \leq 1 \quad (11)$$

These limits correspond for example to the minimum and maximum generator limits if the generator limits are included as part of the input vector. In Section V-B, we will show how including the DC power balance (13) as additional constraint on the input allows to further improve the obtained bounds in (10a).

e) *Regression problems*: Note that the proposed framework can be extended to regression problems. As an example, neural networks can be used to predict the security margin of the power systems (e.g. the damping ratio of the least damped mode in small signal stability analysis). Then, we can solve verification problems of similar structure as (10) and (9) to determine regions in which the neural network prediction is guaranteed to deviate less than a predefined constant from the actual value. At the same time, for a defined region, we can compute the minimum and maximum security margin predicted. This allows us to analyse the robustness of regression neural networks.

IV. IMPROVING TRACTABILITY OF VERIFICATION PROBLEMS

By examining the complexity of the MILPs presented in Section 4, we make two observations. First, from (8a) – (8e) it becomes apparent that the number of required binaries in the MILP is equal to the number of ReLU units in the neural network. Second, the weight matrices \mathbf{W} are dense as the neural network layers are fully connected. As neural networks can potentially have several hidden layers with a large number of hidden neurons, improving the computational tractability of the MILP problems is necessary. To this end, we apply two different methods from [19] and [31]: tightening the bounds on the ReLU output, and sparsifying the weight matrices.

1) *ReLU bound tightening*: In the reformulation of the maximum operator in (8a) and (8c), we introduced upper and lower bounds $\hat{\mathbf{z}}^{\max}$ and $\hat{\mathbf{z}}^{\min}$, respectively. Without any specific knowledge about the bounds, these have to be set to a large value in order not to become binding. A computationally cheap approach to tighten the ReLU bounds is through interval arithmetic (IA) [19]. We propagate the initial bounds on the input $\mathbf{z}_0^{\min} = \mathbf{x}^{\min}$, $\mathbf{z}_0^{\max} = \mathbf{x}^{\max}$ through each layer to obtain individual bounds on each ReLU in layer $k = 1, \dots, K$:

$$\hat{\mathbf{z}}_k^{\max} = \mathbf{W}_{k-1}^+ \hat{\mathbf{z}}_{k-1}^{\max,+} + \mathbf{W}_{k-1}^- \hat{\mathbf{z}}_{k-1}^{\min,+} + \mathbf{b}_{k-1} \quad (12a)$$

$$\hat{\mathbf{z}}_k^{\min} = \mathbf{W}_{k-1}^+ \hat{\mathbf{z}}_{k-1}^{\min,+} + \mathbf{W}_{k-1}^- \hat{\mathbf{z}}_{k-1}^{\max,+} + \mathbf{b}_{k-1} \quad (12b)$$

The max and min-operator are denoted in compact form as: $\mathbf{x}^+ = \max(\mathbf{x}, 0)$ and $\mathbf{x}^- = \min(\mathbf{x}, 0)$. For example, in

our simulation studies, we restrict the input \mathbf{x} to be between $\mathbf{x}^{\max} = \mathbf{0}$ and $\mathbf{x}^{\min} = \mathbf{1}$. Then, the bounds on the first layer evaluate to $\hat{\mathbf{z}}_1^{\max} = \mathbf{W}_0^+ \mathbf{1} + \mathbf{b}_0$ and $\hat{\mathbf{z}}_1^{\min} = \mathbf{W}_0^- \mathbf{1} + \mathbf{b}_0$. The bounds for the remaining layers can be obtained by applying (12) sequentially. Methods to compute tighter bounds also exist, e.g. by solving an LP relaxation while minimizing or maximizing the ReLU bounds [32], but this is out of scope of this paper.

2) Training Neural Networks for Easier Verifiability:

A second possible approach to increase the computational tractability of the verification problems is to (re-)train the neural network with the additional goal to sparsify the weight matrices \mathbf{W} . Here, we rely on an automated gradual pruning algorithm proposed in [31]. Starting from 0% sparsity, where all weight matrices \mathbf{W} are non-zero, a defined share of the weight entries are set to zero. The weight entries selected are those with the smallest absolute magnitude. Subsequently, the neural network is re-trained for the updated sparsified structure. This procedure is repeated until a certain sparsity target is achieved. There are two important observations: First, through sparsification the classification accuracy can decrease, as less degrees of freedom are available during training. Second, larger sparsified networks can achieve better performance than smaller dense networks [31]. As a result, by forming slightly larger neural networks we can maintain the required accuracy while achieving sparsity. As we will see in Section V, sparsification maintains a high classification accuracy and, thus, a significant computational speed-up is achieved when solving the MILPs. As a further benefit of sparsification, the interpretability of the neural network increases, as the only neuron connections kept are the ones that have the strongest contribution to the classification. Computational tractability can further increase by pruning ReLUs during training, i.e. fixing them to be either active or inactive, in order to eliminate the corresponding binaries in the MILP [33]. This approach along with the LP relaxation for bound tightening will be object of our future work.

V. SIMULATION AND RESULTS

A. Simulation Setup

The goal of the following simulation studies is to illustrate the proposed methodology for neural networks which classify operating points as ‘safe’ or ‘unsafe’ with respect to different power system security criteria. The neural network input \mathbf{x} encodes variables such as active generator dispatch and uniquely determines an operating point. The neural network output \mathbf{y} corresponds to the two possible classifications: ‘safe’ or ‘unsafe’ with respect to the specified security criteria.

In the following, we will present four case studies. The first two case studies use a 9-bus and 162-bus system, respectively. The security criteria is feasibility to the N-1 security constrained DC optimality power flow (OPF). The third case study uses a 14-bus system and as security criteria the combined feasibility to the N-1 security constrained AC-OPF and small-signal stability. The fourth case study uses a 162-bus system and as security criteria the feasibility to the N-1 security constrained AC-OPF under uncertainty. The details of the

OPF formulations are provided in the Appendix. Please note that in both the formulation of the N-1 security constrained DC- and AC-OPF we do not consider load shedding, as this should be avoided at all times. Here, the goal is to provide the system operator with a tool to rapidly screen a large number of possible operating points and identify possible critical scenarios. For these critical scenarios only, a more dedicated security constrained OPF could be solved to minimize the cost of load shedding or to redispatch generation units.

To train the neural networks to predict these security criteria, it is necessary to have a dataset of labeled samples that map operating points (neural network inputs \mathbf{x}) to an output security classification $\bar{\mathbf{y}}$. The neural network predicts the output \mathbf{y} which should be as close as possible to the ground truth $\bar{\mathbf{y}}$. We train the neural network as outlined in Section II to achieve satisfactory predictive performance and extract the weights \mathbf{W} and biases \mathbf{b} obtained. Then, based on these, we can formulate the verification problems in (9) and (10) to derive formal guarantees and assess and improve the robustness of these neural networks including the existence of adversarial examples. For a detailed overview of our methodology please refer to the flowchart in Figure 2.

The created dataset can include both historical and simulated operating points. Absent historical data in this paper, we created simulated data for our studies. We will detail the dataset creation for each case study in the corresponding subsection. To facilitate a more efficient neural network training procedure, we normalize each entry of the input vector \mathbf{x} to be between 0 and 1 using the upper and lower bounds of the variables (we do that for all samples $\mathbf{x} \in \mathbf{S}$). Empirically, this has been shown to improve classifier performance.

After creating the training dataset, we export it to TensorFlow [25] for the neural network training. We split the dataset into a training set and a test set. We choose to use 85% of samples for training and 15% for testing. During training, we minimize the cross-entropy loss function (5) using the Adam optimizer with stochastic gradient descent [25], and use the default options in TensorFlow for training. In the cases where we need to enforce a certain sparsity of the weight matrices we re-train with the same objective function (5), but during training we reduce the number of non-zero weights until a certain level of sparsity is achieved, as explained in Section IV-2. To allow for the neural network verification and the identification of adversarial examples, after the neural network training we export the weight matrices \mathbf{W} and biases \mathbf{b} in YALMIP, formulate the MILPs, and solve them with Gurobi. If not noted otherwise, we solve all MILPs to zero optimality gap, and as a result obtain the globally optimal solution. All simulations are carried out on a laptop with processor Intel(R) Core(TM) i7-7820HQ CPU @ 2.90 GHz and 32GB RAM.

B. Neural Network Verification for the IEEE 9-bus system

1) *Test Case Setup:* For the first case study, we consider a modified IEEE 9-bus system shown in Fig. 3. This system has three generators and one wind farm. We assume that the load at bus 7 is uncertain and can vary between 0 MW and 200 MW. The wind farm output is uncertain as well and

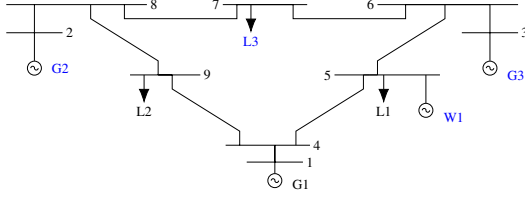


Fig. 3. Modified IEEE 9-bus system with wind farm $W1$ at bus 5, and uncertain load $L3$ at bus 7. Generation and loading marked in blue correspond to the input vector \mathbf{x} of the neural network.

can vary between 0 MW and 300 MW. The line limits are increased by 25% to accommodate the added wind generation and increased loading. Remaining parameters are defined according to [34]. As mentioned in the previous subsection, in this case we employ the N-1 security for the DC power flow as a security criterion. If the SC-DC-OPF results to infeasibility, then the sample is ‘unsafe’. We consider the single outage of each of the six transmission lines and use MATPOWER to run the DC power flows for each contingency fixing the generation and load to the defined operating point. We evaluate the violation of active generator and active line limits to determine whether an operating point is ‘safe’ or ‘unsafe’. We create a dataset of 10’000 possible operating points with input vector $\mathbf{x} = [P_{G2} P_{G3} P_{L3} P_{W1}]$, which are classified $\mathbf{y} = \{y_1 = \text{safe}, y_2 = \text{unsafe}\}$. Note that the first generator P_{G1} is the slack bus and its power output is uniquely determined through the dispatch of the other units; therefore, it is not considered as input to the neural network. The possible operating points are sampled using Latin hypercube sampling [35] which ensures that the distance between each sample is maximized.

To compute the ground-truth (and thus be able to assess the neural network performance) we use a custom implementation of the N-1 preventive security constrained DC-OPF (SC-DC-OPF) in YALMIP using pre-built functions of MATPOWER to compute the bus and line admittance matrix for the intact and outaged system state. We model the the uncertain loads and generators as optimization variables. With ground truth, we refer to the region around an infeasible sample for which we can guarantee that no input exists which is feasible to the SC-DC-OPF problem. We can compute this by minimizing the distance from an infeasible sample to a feasible operation point. For the detailed mathematical formulation of the N-1 security constrained DC-OPF and the ground truth evaluation please refer to Appendix A-B.

2) *Neural Network Training*: Using the created dataset, we train two neural networks which only differ with respect to the enforced sparsity of the weight matrices. Employing loss function (5), the first neural network is trained with dense weight matrices. Using this as a starting point for the retraining, on the second neural network we enforce 80% sparsity, i.e. only 20% of the entries are allowed to be non-zero. In both cases, the neural network architecture comprises three hidden layers with 50 neurons each, as this allows us to achieve high accuracy without over-fitting. For comparison, if we would only train a single-layer neural network with 150 neurons, the maximum test set classification accuracy is only 90.7%, highlighting the need for a multi-layer architecture.

TABLE I
CONFUSION MATRICES FOR IEEE 9-BUS TEST CASE

Neural network <u>without</u> sparsified weight matrices			
Test samples: 1500	Predicted: Safe	Unsafe	Accuracy
True: Safe (326)	311	15	
True: Unsafe (1174)	11	1163	
			98.3%
Neural network <u>with</u> 80% sparsified weight matrices			
Test samples: 1500	Predicted: Safe	Unsafe	Accuracy
True: Safe (326)	308	18	
True: Unsafe (1174)	15	1159	
			97.8%

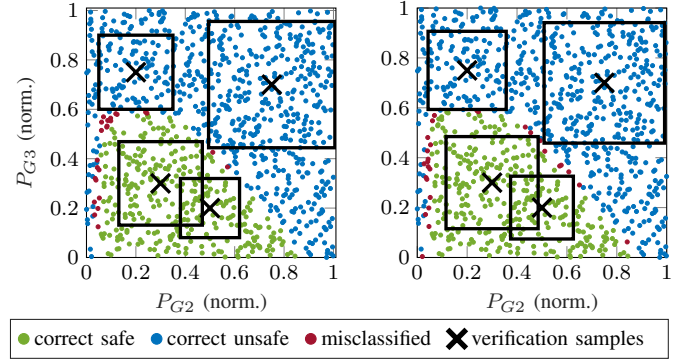


Fig. 4. Regions around four verification samples in which the classification is guaranteed not to change. The left figure uses the neural network without sparse weight matrices, and the right figure uses the neural network with imposed 80% sparsity on the weight matrices. For visualization purposes, the load P_{L3} and wind level P_{W1} are fixed to 40% and 1’000 new samples are created and classified according to the respective neural network.

Our three-layer neural network without sparsity has a classification accuracy of 98.3% on the test set and 98.4% on the training set. This illustrates that neural networks can predict with high accuracy whether an operating point is safe or unsafe. If we allow only 20% of the weight matrix entries to be non-zero and re-train the three-layer neural network, the classification accuracy evaluates to 97.8% on the test set and 97.3% on the training set. The corresponding confusion matrices for both neural networks for the test set are shown in Table I. As we see, the sparsification of the neural network leads to a small reduction of 0.5% in classification accuracy, since less degrees of freedom are available during training. The benefits from sparsification are, however, substantially more significant. As we will see later in our analysis, sparsification substantially increases the computational tractability of verification problems, and leads to an increase in interpretability of the neural network.

3) *Visualization of Verification Samples*: To be able to visualize the input regions for which we can guarantee that the trained neural networks will not change their classification, we shall reduce the input dimensionality. For this reason, we will study here only a single uncertainty realization, where we assume that both the load P_{L3} and wind power output P_{W1} amount to 40% of their normalized output (note that our results apply to all possible inputs). The resulting input space is two-dimensional and includes generators P_{G2} and

P_{G3} . For visualization purposes only, we take 1'000 new samples from the reduced two-dimensional space and classify them with the neural networks as safe or unsafe. In addition, we compute their true classification using MATPOWER. The resulting classified input spaces are shown in Fig. 4 with the left figure corresponding to the neural network with full weight matrices, and the right figure to the sparsified weight matrices. We can observe that misclassifications occur at the security boundary. This is to be expected as the sampling for this visualization is significantly more dense than the one used for training. What is important here to observe though, is that even if the neural network on the right contains only 20% of the original number of non-zero weights, the neural networks have visually comparable performance in terms of classification accuracy. As a next step, we solve the verification problem described in (10). In that, for several samples \mathbf{x}_{ref} we compute the minimum perturbation ϵ that changes the classification. We visualize the obtained regions using black boxes around the verification samples in Fig. 4. The average solving time in Gurobi is 5.5 s for the non-sparsified and 0.3 s for the sparsified neural network. We see that by sparsifying the weights we achieve a $20\times$ computational speed-up. Solving the MILP to zero optimality gap, we are guaranteed that within the region defined by the black boxes around the verification samples no input exists which can change the classification.

4) *Provable Guarantees:* In the following, we want to provide several provable guarantees of behaviour for both neural networks while considering the entire input space. For this purpose, we solve verification problems of the form (10). Since the neural network input \mathbf{x} is normalized between $\mathbf{0} \leq \mathbf{x} \leq \mathbf{1}$ based on the bounds of the input variables, we include the corresponding constraint (11) in the verification problem. Similarly, no input to the neural network would violate the limits of the slack bus generator P_{G1} . These inputs could be directly classified as unsafe. Even if P_{G1} is not part of the input, its limits can be defined based on the input \mathbf{x} through the DC power balance, as shown in (13):

$$P_{G1}^{\min} \leq P_{L1} + P_{L2} + P_{L3} - P_{G2} - P_{G3} - P_{W1} \leq P_{G1}^{\max}, \quad (13)$$

where $\mathbf{x} = [P_{G2}, P_{G3}, P_{L3}, P_{W1}]^T$. The DC power balance ensures that the system generation equals the system loading. As a result, the slack bus has to compensate the difference between system loading and generation. In this section, we show how additional a-priori qualifications of the input such as (13), affect the size of the regions in which we can guarantee a certain classification.

In Table II, we compare the obtained bounds to the ground truth provided by the SC-DCOPF and report the computational time for three properties, when we do not include and when we do include the power balance constraint. The second reported computational time uses the tightened ReLU bounds (12) with interval arithmetic (IA) (see Section IV-1 for details on the method). The first property is the size of the largest region around the operating point $\mathbf{x}_{\text{ref}} = \mathbf{1}$ (1-vector) with the same guaranteed classification. This operating point corresponds to the maximum generation and loading of the system. For this input, the classification of the neural networks is known to be

TABLE II
VERIFICATION OF NEURAL NETWORK PROPERTIES FOR 9-BUS SYSTEM

	w/o power balance		with power balance	
	ϵ	Sol. Time (s)	ϵ	Sol. Time (s)
Property 1: $\forall \mathbf{x} \in [\mathbf{0}, \mathbf{1}] : \mathbf{x} - \mathbf{1} _{\infty} \leq \epsilon \rightarrow \text{Classification: insecure}$				
NN (w/o sparsity)	50.7%	2.7 IA: 1.1	54.4%	1.5 IA: 1.3
NN (80% sparsity)	48.7%	0.3 IA: 0.2	54.4%	0.3 IA: 0.2
SC-DC-OPF	53.7%	-	53.7%	-
Property 2: $\forall \mathbf{x} \in [\mathbf{0}, \mathbf{1}] : \mathbf{x} - \mathbf{0} _{\infty} \leq \epsilon \rightarrow \text{Classification: secure}$				
NN (w/o sparsity)	29.2%	501.9 IA: 400.7	29.3%	473.0 IA: 817.9
NN (80% sparsity)	32.7%	3.3 IA: 1.1	32.7%	2.1 IA: 1.6
SC-DC-OPF ¹	31.7%	-	31.7%	-
Property 3: $\exists \mathbf{x} \in [\mathbf{0}, \mathbf{1}] : \mathbf{P}_{W2} - \mathbf{0} _{\infty} \leq \epsilon \rightarrow \text{Classification: secure}$				
NN (w/o sparsity)	97.4%	12.2 IA: 11.4	90.3%	6.5 IA: 3.7
NN (80% sparsity)	99.1%	0.4 IA: 0.3	89.0%	0.4 IA: 0.3
SC-DC-OPF	92.5%	-	92.5%	-

¹ The SC-DC-OPF can only provide infeasibility certificates.

We compute this by re-sampling a very large number of samples.

unsafe. We observe for both neural networks that when the power balance constraint is not included in the verification problem, the input region guaranteed to be classified as unsafe is smaller than the ground truth of 53.7% (provided by the SC-DC-OPF). By including the power balance in the verification problem, the input region classified as unsafe is significantly enlarged and more closely matches the target of 53.7%.

For the second property, we consider the region around the operating point where all generating units and the load are at their lower bound, i.e., $\mathbf{x}_{\text{ref}} = \mathbf{0}$. This point corresponds to a low loading of the system and is therefore secure (i.e. no line overloadings, etc.). We solve the verification problem minimizing the distance from $\mathbf{x}_{\text{ref}} = \mathbf{0}$ leading to an insecure classification. In this way, we obtain the maximum region described as hyper-cube around the zero vector in which the classification is guaranteed to be always ‘safe’. We can observe in this case that the neural network without sparsity slightly underestimates while the neural network with 80% sparsity slightly overestimates the safe region compared to the ground truth (31.7%). Sparsification allows a $150\times$ – $500\times$ computational speed-up (two orders of magnitude). For this property, including the power balance constraint does not change the obtained bounds.

The third property we analyze is the maximum wind infeed for which the neural network identifies a secure dispatch. To this end, we solve the verification problem by maximizing the wind infeed P_{W1} in the objective function while maintaining a secure classification. The true maximum wind infeed obtained by solving this problem directly with the SC-DC-OPF is 92.5%, i.e. 277.5 MW can be accommodated. If we do not enforce constraint (13), then we observe that the obtained bounds 97.4% and 99.12% are not tight. This happens because the obtained solutions \mathbf{x} from (10) keep generation and loading to zero and only maximize wind power; this violates the generator bounds on the slack bus, as it has to absorb all the generated wind power. Enforcing the DC power balance (13) in the verification problem allows a more accurate representation of the true admissible inputs, and we obtain significantly improved bounds of 90.3% and 89.0%.

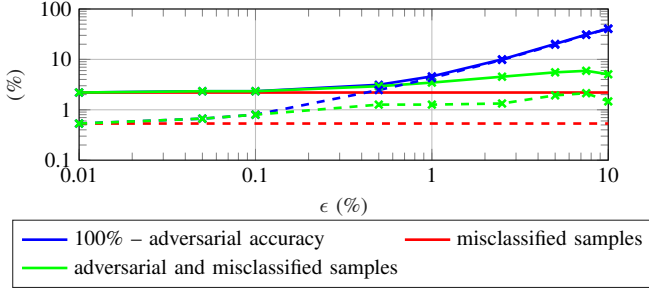


Fig. 5. Adversarial accuracy, and share of adversarial and misclassified samples are shown for the test data set of the 9 bus test case and different levels of input perturbation ϵ . The adversarial accuracy refers to the share of correctly classified samples, for which no input perturbation exists which changes the classification within distance ϵ of that sample (by solving (9) for each sample). Please note that both axes are logarithmic, and 100% minus the adversarial accuracy is shown, i.e. the share of samples which are not adversarially robust. Out of these, to determine whether an adversarial example has been identified, the ground-truth classification is computed. The dashed lines show the performance of the re-trained neural network when adding 1'152 identified adversarial examples from the training dataset to the updated training dataset. It can be observed that both prediction accuracy and adversarial robustness, i.e. the share of adversarial examples, are significantly improved.

For all three properties, we observe that both interval arithmetic (IA) and weight sparsification improve computational tractability, in most instances achieving the lowest solver times when used combined. From the two, weight sparsification has the most substantial effect as it allows to reduce the number of binaries in the resulting MILP problem.

5) *Adversarial robustness and re-training*: In Fig. 5, the adversarial robustness and the share of adversarial and misclassified samples are shown for the test dataset and different magnitudes of input perturbation ϵ in the range from 0.01% to 10%, following the methodology described in Fig. 2. The full lines are computed using the sparsified neural network with performance described in the confusion matrix in Table I. The sparsified neural network has a predictive accuracy of 97.8% on the training dataset, i.e. 2.2% of test samples are misclassified. It can be observed that for small input perturbations ϵ the share of adversarial and misclassified samples increases to 3.5%, i.e. for an additional 1.3% of test samples the identified adversarial input leads to a wrong misclassification. Increasing the input perturbation ϵ to 10%, these shares increase to 5.1% and 2.9%, respectively. This indicates that parts of the security boundary are not correctly represented by the neural network.

To improve performance, we run the same methodology for the training dataset of 8'500 samples, and include all identified adversarial examples of the training dataset (a total of 1'152 samples) in the updated training dataset and re-train the neural network. Note that we only use the training dataset for this step as the unseen test dataset is used to evaluate the neural network performance. The resulting performance is depicted with dashed lines in Fig. 5. Re-training the neural network has two benefits in this test case: First, it improves the predictive accuracy on the unseen test data from 97.8% to 99.5%, i.e. only 0.5% of test samples are misclassified. Second, the share of identified adversarial samples is reduced, for $\epsilon = 1\%$ from 1.3% to 0.7% and for $\epsilon = 1\%$ from 2.9% to 0.9%, showing improved neural network robustness. At the same time, it can be observed that for perturbations larger than $\epsilon = 1\%$, the

TABLE III
CONFUSION MATRIX FOR IEEE 162-BUS TEST CASE (WITH SPARSITY)

Test samples: 3000	Predicted: Safe	Unsafe	Accuracy
True: Safe (1507)	1501	6	
True: Unsafe (1493)	11	1482	
Accuracy			99.4%

adversarial accuracy is similar between both networks. This shows that in this case for a large amount of samples an adversarial input can be identified which correctly leads to a different classification, i.e. the adversarial input moves the operating point across the true security boundary and is not an adversarial example leading to a false misclassification. The neural network robustness could be further improved by repeating this procedure for additional iterations.

C. Scalability for IEEE 162-bus system

1) *Test Case Setup*: For the second case study, we consider the IEEE 162-bus system with parameters taken from [34]. We add five uncertain wind generators located at buses {3, 20, 25, 80, 95} with a rated power of 500 MW each. As security criterion, we consider again N-1 security based on DC power flow, considering the outage of 24 critical lines: {22, ..., 27, 144, ..., 151, 272, ..., 279}. The input vector is defined as $\mathbf{x} = [P_{G1-G12} P_{W1-W5}]^T$. To construct the dataset for the neural network training, we first apply a bound tightening of the generator active power bounds, i.e. we maximize and minimize the corresponding bound considering the N-1 SC-DC-OPF constraints. The tightened bounds exclude regions in which the SC-DC-OPF is guaranteed to be infeasible and therefore allows to decrease the input space. In the remaining input space, we draw 10'000 samples using Latin hypercube sampling and classify them as safe or unsafe. As usual in power system problems, the 'safe' class is substantially smaller than the 'unsafe' class, since the true safe region is only a small subset of the eligible input space. To mitigate the dataset imbalance, we compute the closest feasible dispatch for each of the infeasible samples by solving an SC-DC-OPF problem and using the ∞ -norm as distance metric. As a result, we obtain a balanced dataset of approximately 20'000 samples.

2) *Neural Network Training*: We choose a neural network architecture with 4 layers of 100 ReLU units at each layer, as the input dimension increases from 4 to 17 compared to the 9-bus system. We train again two neural networks: one without enforcing sparsity and a second with 80% sparsity. The trained four-layer network without sparsity has a classification accuracy 99.3% on the test set and 99.7% on the training set. For the sparsified network, these metrics evaluate to 99.7% for the test set and 99.4% for the training set. The confusion matrix of the sparsified network is shown in Table III and it can be observed that the neural network has high accuracy in classifying both safe and unsafe operating points.

3) *Provable Guarantees*: Assuming the given load profile, the property of interest is to determine the minimum distance from zero generation $\mathbf{x}_{\text{ref}} = \mathbf{0}$ to a feasible ('safe') solution.

TABLE IV
VERIFICATION OF NEURAL NETWORK PROPERTIES FOR 162-BUS SYSTEM

	w/o power balance		with power balance	
	ϵ	Sol. Time (s)	ϵ	Sol. Time (s)
Property 1: $\forall x \in [0, 1] : \ x - 0\ _\infty \leq \epsilon \rightarrow \text{Classification: insecure}$				
NN (w/o sparsity)	-	> 50 min	-	> 50 min
NN (80% sparsity)	59.4%	560 IA: 1217	65.5%	22 IA: 30
SC-DC-OPF	66.2%	-	66.2%	-

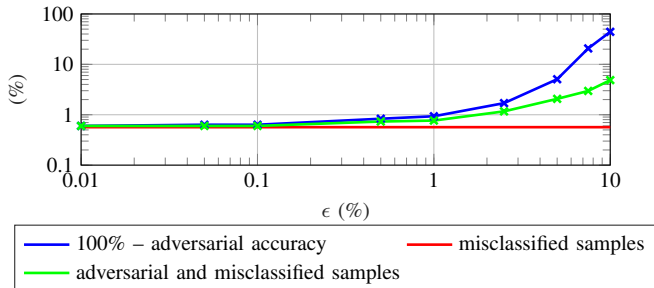


Fig. 6. Adversarial accuracy, and share of adversarial and misclassified samples are shown for the test data set of the 162 bus test case and different levels of input perturbation ϵ . The adversarial accuracy refers to the share of correctly classified samples, for which no input perturbation exists which changes the classification within distance ϵ of that sample (by solving (9) for each sample). Please note that both axes are logarithmic, and 100% minus the adversarial accuracy is shown, i.e. the share of samples which are not adversarially robust. Out of these, to determine whether an adversarial example has been identified, the ground-truth classification is computed.

In Table IV, we compare the result of the verification problem with the ground truth obtained from solving the SC-DC-OPF problem. We can see that the bound of 59.4% we obtain without including the DC power balance is not tight compared to the ground truth of 66.2%. Including the DC power balance increases the bound to 65.2% which is reasonably close, indicating satisfactory performance with respect to this property. This confirms the findings for the 9-bus (cmp. Tables II and IV) showing that including the additional power balance constraint in the verification problems leads to bounds on ϵ closer to the ground-truth. Regarding computational time, we can observe that for this larger neural network, the sparsification of the weight matrices becomes a requirement to achieve tractability. For both cases with and without including the DC power balance, the MILP solver did not identify a solution after 50 minutes for the non-sparsified network.

4) *Adversarial Robustness*: In Fig. 6, the adversarial accuracy, and share of adversarial and misclassified samples are computed for the test data set and different levels of input perturbation ϵ . It can be noted that for small input perturbations $\epsilon \leq 1\%$ the adversarial accuracy is above 99% (i.e. for $> 99\%$ of test samples no input exists which can change the classification), indicating neural network robustness. For larger input perturbations, the adversarial accuracy decreases as system input perturbations can move the operating point across system security boundaries. The number of identified adversarial examples increases as well. For a large input perturbation $\epsilon = 10\%$ for 129 (4.3%) test samples an adversarial input exists which falsely changes the classification. This indicates that the security boundary estimation of the neural network is inaccurate in some parts of the high-dimensional

TABLE V
CONFUSION MATRIX FOR IEEE 14-BUS TEST CASE (WITH SPARSITY)

Test samples: 1500	Predicted: Safe	Unsafe	Accuracy
True (15): Safe	12	3	
True (1485): Unsafe	0	1485	
Accuracy			99.8%

input space. As for the 9 bus test case, a re-training of the neural network by including identified adversarial examples in the training data set could improve performance.

D. N-1 Security and Small Signal Stability

Security classifiers using neural networks can screen operating points for security assessment several orders of magnitudes faster than conventional methods [15]. There is, however, a need to obtain guarantees for the behaviour of these classifiers. We show in the following how our framework allows us to analyse the performance of a classifier for both N-1 security and small signal stability.

1) *Dataset Creation*: For the third case study, we consider the IEEE 14-bus system [34] and create a dataset of operating points which are classified according to their feasibility to *both* the N-1 SC-AC-OPF problem and small signal stability. We consider the outages of all lines except lines 7-8 and 6-13, as the 14 bus test case is not N-1 secure for these two line outages. Furthermore, we assume that if an outage occurs, the apparent branch flow limits are increased by 12.5% resembling an emergency overloading capability. We use a simple brute-force sampling strategy to create a dataset of 10'000 equally spaced points in the four input dimensions $\mathbf{x} = [P_{G2-G5}]$. The generator automatic voltage regulators (AVRs) are fixed to the set-points defined in [34]. For each of these operating points, we run an AC power flow and the small-signal stability analysis for each contingency. Operating points which satisfy operational constraints and the eigenvalues of the linearized system matrix have negative real parts for all contingencies are classified as 'safe' and 'unsafe' otherwise. Note that, similar to the 162-bus case, the dataset is unbalanced as only 1.36% of the overall created samples are feasible. A method to create balanced datasets for *both* static and dynamic security is object of future work. For the full mathematical details on the N-1 SC-AC-OPF formulation and the small-signal stability constraints, please refer to Appendices A-C and A-D, respectively.

2) *Neural Network Training and Performance*: We choose a three-layer neural network with 50 neurons for each layer, as in Section V-B. Based on the analysis of the previous case studies, here we only train an 80% sparsified network. The network achieves the same accuracy of 99.8% both on the training and on the test set. The confusion matrix on the test set is shown in Table V. Note that here the accuracy does not carry sufficient information as a metric of the neural network performance, as simply classifying all samples as unsafe would lead to a test set accuracy of 99.0% due to the unbalanced classes. Besides the use of supplementary metrics, such as specificity and recall (see e.g. [15]), obtaining provable

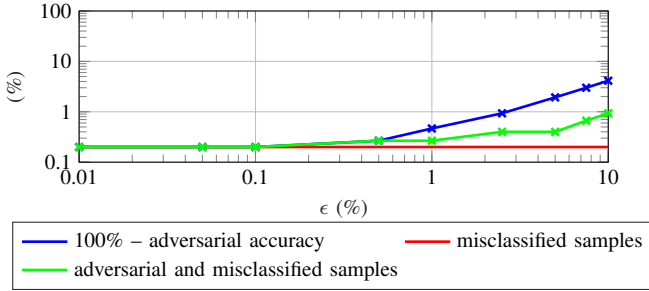


Fig. 7. Adversarial accuracy, and share of adversarial and misclassified samples are shown for the test data set of the 14 bus test case and different levels of input perturbation ϵ . The adversarial accuracy refers to the share of correctly classified samples, for which no input perturbation exists which changes the classification within distance ϵ of that sample (by solving (9) for each sample). Please note that both axes are logarithmic, and 100% minus the adversarial accuracy is shown, i.e. the share of samples which are not adversarially robust. Out of these, to determine whether an adversarial example has been identified, the ground-truth classification is computed.

guarantees of the neural network behavior becomes, therefore, of high importance.

One of the main motivations for data-driven approaches are their computational speed-up compared to conventional methods. To assess the computational complexity, we compare the computational time required for evaluating 1'000 randomly sampled operating points using AC power flows and the small-signal model to using the trained neural network to predict power system security. We find that, on average, computing the AC power flows and small-signal model for all contingencies takes 0.22 s, and evaluating the neural network 7×10^{-5} s, translating to a computational speed up of factor 1000 (three orders of magnitude). Similar computational speed-ups are reported in [15], [16] for deep learning applications to system security assessment. Note that the works in [15], [16] do not provide performance guarantees and do not examine robustness. Contrary, in the following, we provide formal guarantees for the performance of the security classifier by analysing the adversarial accuracy and identifying adversarial examples.

3) *Evaluating Adversarial Accuracy*: Adversarial accuracy identifies the number of samples whose classification changes from the ground-truth classification if we perform a perturbation to their input. Assuming that points in the neighborhood of a given sample should share the same classification, e.g. points around a 'safe' sample would likely be safe, a possible classification change indicates that we are either very close to the security boundary or we might have discovered an adversarial example. Carrying out this procedure for our whole test dataset, we would expect that in most cases the classification in the vicinity of the sample will not change (except for the comparably small number of samples that is close to the security boundary). In Fig. 7 the adversarial accuracy is depicted for the sparsified neural network. It can be observed that for small perturbations, i.e. $\epsilon \leq 1\%$, the adversarial accuracy stays well above 99%; that is, for 99% of test samples no input exists within distance of ϵ which changes the classification. Only if large perturbations to the input are allowed (i.e. $\epsilon \geq 1\%$) the adversarial accuracy decreases. This shows that the classification of our neural network is adversarially robust in most instances. Note that the adversarial

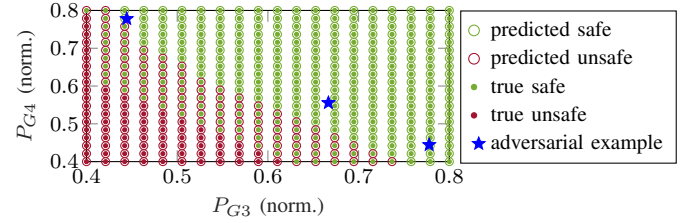


Fig. 8. Classification of 400 new equally spaced samples for the IEEE 14-bus system. For 2-D visualization purposes, the active power of P_{G2} and P_{G5} are fixed to their maximum output. Blue stars mark identified adversarial examples. For these samples a small input perturbation exists which falsely changes the classification. The reason is the inaccurate prediction of the system security boundary.

accuracy only makes a statement regarding the classification by the neural network and is not related to the ground truth classification. In a subsequent step, as we show in the next paragraph, we need to assess whether identified samples are in fact misclassified.

4) *Identifying Adversarial Examples*: Having identified regions where adversarial examples may exist, we can now proceed with determining if they truly exist. The resulting share of identified adversarial examples is shown in Fig. 7. Focusing on the test samples which are not adversarially robust for small ϵ , i.e. $\epsilon \leq 1\%$, we identify an adversarial example for the sample $\mathbf{x} = [1.0 \ 1.0 \ 0.4444 \ 0.7778]$ with classification 'safe'. Modifying this input only by $\epsilon = 0.5\%$, we identify the adversarial input $\mathbf{x}_{adv} = [0.9950 \ 0.9950 \ 0.4394 \ 0.7728]$ which falsely changes the classification to 'unsafe', i.e. $y_{adv,1} > y_{adv,2}$. Allowing an input modification of magnitude $\epsilon = 1\%$, we identify two additional adversarial examples for inputs $\mathbf{x} = [1.0 \ 1.0 \ 0.6667 \ 0.5556]$ and $\mathbf{x} = [1.0 \ 1.0 \ 0.7778 \ 0.4444]$ with classification 'safe', respectively. These have the corresponding adversarial inputs $\mathbf{x}_{adv} = [0.9750 \ 0.9750 \ 0.6417 \ 0.5306]$ and $\mathbf{x}_{adv} = [0.9750 \ 0.9750 \ 0.7528 \ 0.4194]$ which falsely change the classification to 'unsafe', respectively.

For illustrative purposes in Fig. 8, we re-sample 400 equally spaced samples and compute both the neural network prediction and ground truth. In Fig. 8 the location of the adversarial examples is marked by a star. We can observe that the neural network boundary prediction is not precise and as a result, the adversarial inputs get falsely classified as unsafe. This highlights the additional benefit of the presented framework to identify adversarial examples, and subsequently regions in which additional detailed sampling for re-training the classifier is necessary.

E. N-1 Security and Uncertainty

For the fourth case study, to further demonstrate scalability of our methodology, we use the IEEE 162 bus test case with parameters defined in [36], and train a neural network to predict power system security with respect to the N-1 security constrained AC-OPF under uncertainty. Compared to the previous 14 bus test case, we assume that the voltage set-points of generators can vary within their defined limits (i.e. they are part of the input vector \mathbf{x}), and we consider both uncertain injections in power generation and demand. For the N-1 security assessment, we consider the possible outages of lines $\{6, 8, 24, 50, 128\}$, assuming the same parameters

TABLE VI
CONFUSION MATRIX FOR IEEE 162-BUS TEST CASE (N-1 SECURITY
AC-OPF AND UNCERTAINTY)

Test samples: 18204	Predicted: Safe	Unsafe	Accuracy
True: Safe (4260)	3148	1112	
True: Unsafe (13944)	609	13335	
Accuracy			90.5%

for the outaged system state as for the intact system state. Furthermore, we place 3 wind farms with rated power of 500 MW and consider 3 uncertain loads with $\pm 50\%$ variability, i.e., a total of 6 uncertain power injections, at buses $\{60, 90, 145, 3, 8, 52\}$. For all uncertain injections, we assume a power factor $\cos \phi = 1$.

1) *Dataset creation:* As the resulting input dimension of the considered test case is high (29 inputs) and large parts of the input space correspond to infeasible operating points, a sampling strategy based on Latin hypercube sampling from the entire input space is not successful at recovering feasible samples. To create the dataset, we rely on an efficient dataset creation method proposed in [37]. The dataset creation method consists of several steps. First, all the upper and lower bounds in the AC-OPF problem and on the input variables in \mathbf{x} are tightened using bound tightening algorithms in [38] and [39]. Second, relying on infeasibility certificates based on convex relaxations of the AC-OPF, large parts of the input space can be characterized as infeasible, i.e. ‘unsafe’ with respect to the power system security criteria. These infeasibility certificates take the form of hyperplanes and the remaining unclassified space can be described as a convex polyhedron $\mathbf{Ax} \leq \mathbf{b}$. Third, a large number of samples (here 10’000) are sampled uniformly from inside the convex polyhedron and their classification is computed. For infeasible samples, the closest feasible sample is determined to characterize the security boundary. Finally, a Gaussian distribution is fitted to the feasible boundary samples, and an additional large number of samples (here 100’000) are sampled from this distribution and their feasibility to the N-1 security constrained AC-OPF is assessed. This methodology results to a database of 121’358 samples out of which 23.2% correspond to ‘safe’ operating points. Note that computing the infeasibility certificates allows to reduce the considered input volume from a normalized volume of 1 (i.e. all bounds on \mathbf{x} are between 0 and 1) to a volume of 6×10^{-10} . This highlights the need for advanced methods for dataset creation, as direct sampling strategies are not able to produce balanced datasets balanced of ‘safe’ and ‘unsafe’ operating points. For more details on the dataset creation method, for brevity, the reader is referred to [37].

2) *Neural Network Training and Adversarial Robustness:* We train a neural network with 3 hidden layers of 100 neurons each and enforce a weight sparsity of 80%. The neural network has an accuracy of 91.3% on the training dataset and 90.5% on the test dataset. The confusion matrix for the test dataset is shown in Table VI. Similar to previous test cases, we evaluate the adversarial accuracy in Fig. 9. We find that the neural network is not adversarially robust, already for an input modification of $\epsilon = 0.1\%$ for 4.2% of test samples an adversarial

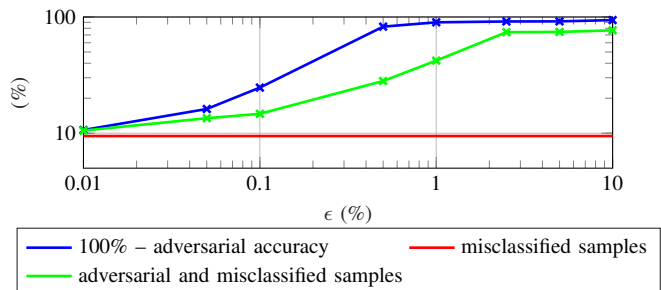


Fig. 9. Adversarial accuracy, and share of adversarial and misclassified samples are shown for the test data set of the 162 bus test case (N-1 security and uncertainty) and different levels of input perturbation ϵ . The adversarial accuracy refers to the share of correctly classified samples, for which no input perturbation exists which changes the classification within distance ϵ of that sample (by solving (9) for each sample). Please note that both axes are logarithmic, and 100% minus the adversarial accuracy is shown, i.e. the share of samples which are not adversarially robust. Out of these, to determine whether an adversarial example has been identified, the ground-truth classification is computed.

example is identified (in addition to the initially misclassified 9.5% of test samples). For an input modification of $\epsilon = 1\%$ this value increases to 31.6%. This systematic process to identify adversarial examples proposed in this paper allows us to obtain additional insights about the quality of the training database. Assessing the adversarial accuracy (i.e. the number of samples that change classification within a distance ϵ) versus the actual adversarial examples, we find that the change in classification in this case often occurs because the operating points have been moved across the true system security boundary. This indicates that many samples are placed close to the correctly predicted true system boundary in the high-dimensional space. Using high-performance computing, an additional detailed re-sampling of the system security boundary or re-training by including adversarial examples as shown in Section V-B5 could improve neural network robustness.

VI. CONCLUSION

Neural networks in power system applications have so far been treated as a black box; this has become a major barrier towards their application in practice. This is the first work to present a rigorous framework for neural network verification in power systems and to obtain provable performance guarantees. To this end, we formulate verification problems as mixed-integer linear programs and train neural networks to be computationally easier verifiable. We provably determine the range of inputs that are classified as safe or unsafe, and systematically identify adversarial examples, i.e. slight modifications in the input that lead to a mis-classification by neural networks. This enables power system operators to understand and anticipate the neural network behavior, building trust in them, and remove a major barrier toward their adoption in power systems. We verify properties of a security classifier for an IEEE 9-bus system, improve its robustness and demonstrate its scalability for a 162-bus system, highlighting the need for sparsification of neural network weights. Finally, we further identify adversarial examples and evaluate the adversarial accuracy of neural networks trained to assess N-1 security under uncertainty and small-signal stability.

APPENDIX A

OPTIMAL POWER FLOW (OPF) FORMULATIONS

In the following, for completeness, we provide a brief overview of the DC and AC optimal power flow formulations including N-1 security and small-signal stability. For more details please refer to [40]–[42].

A. Preliminaries

We consider a power grid which consists of buses (denoted with the set \mathcal{N}) and transmission lines (denoted with the set \mathcal{L}). The transmission lines connect one bus $i \in \mathcal{N}$ to another bus $j \in \mathcal{N}$, i.e., $(i, j) \in \mathcal{L}$. For the AC-OPF formulation, we consider the following variables for each bus: The voltage magnitudes V , the voltage angles θ , the active power generation P_G , the reactive power generation Q_G , the active wind power generation P_W , and the active and reactive power demand P_L and Q_L . Each of these vectors have the size $n_b \times 1$, where n_b is the number of buses in the set \mathcal{N} . Note that during operation, for one specific instance, the active wind power generation and active and reactive power demand are assumed to be fixed (and the curtailment of wind generators and load shedding are to be avoided at all times). To comply with the N-1 security criterion, we consider the possible single outage of a set of lines, which we denote with the set \mathcal{C} . The first element of this set corresponds to the intact system state, denoted with '0'. The superscript 'c' denotes the variables corresponding to the intact and outaged system states.

B. N-1 Security-Constrained DC-OPF

In Section V-B and Section V-C, we create datasets of operating points classified according to their feasibility to the N-1 security constrained DC-OPF. The DC-OPF approximation neglects reactive power and active power losses, and assumes that the voltage magnitudes V are fixed, for a detailed treatment please refer to [41]. Considering a set \mathcal{C} of possible single line outages, the preventive N-1 security constrained DC-OPF problem can then be formulated as follows:

$$\min_{P_G^c, P_L, P_W, \theta^c} f(P_G^0) \quad (14)$$

$$\text{s.t. } P_G^c + P_W - P_L = B_{\text{bus}}^c \theta^c \quad \forall c \in \mathcal{C} \quad (15)$$

$$P_{\text{line}}^{\min, c} \leq B_{\text{line}}^c \theta^c \leq P_{\text{line}}^{\max, c} \quad \forall c \in \mathcal{C} \quad (16)$$

$$P_G^{\min} \leq P_G^c \leq P_G^{\max} \quad \forall c \in \mathcal{C} \quad (17)$$

$$P_W^{\min} \leq P_W \leq P_W^{\max} \quad (18)$$

$$P_L^{\min} \leq P_L \leq P_L^{\max} \quad (19)$$

$$P_G^{\text{wosb}, 0} = P_G^{\text{wosb}, c} \quad \forall c \in \mathcal{C} \quad (20)$$

The objective function f minimizes e.g. the generation cost of the intact system state. The nodal power balance in (15) has to be satisfied for the intact and outaged system states. The bus and line admittance matrices are denoted with B_{bus} and B_{line} , respectively. Upper and lower limits are enforced for the active power line flows, generation, wind power, and load demands in (16), (17), (18) and (19), respectively. The constraint in (20) enforces preventive control action of generators. The superscript 'wosb' denotes all generators except the

slack bus generator. The independent variables characterizing an operating point are $\mathbf{x} := [P_G^{\text{wosb}, 0}, P_W, P_L]$. To create datasets, we compute the classification for each operation point by first running DC power flows to determine θ^c and the slack bus generator dispatch $P_G^{\text{slack}, c}$ for the intact and each outaged system states. The superscript 'slack' denotes the slack bus. Then, we check satisfaction of the constraints on active generator power (17) and active line flows (16). In operations, it is usually assumed that both the wind power and loading are fixed, i.e. $P_W^{\max} = P_W^{\min}$ and $P_L^{\max} = P_L^{\min}$. Here, we model them as variables to be able to compute the ground-truth for the region around an infeasible sample $\mathbf{x}^{\text{infeas}}$ in which no feasible sample exist. To this end, we solve the following optimization problem computing the minimum distance from the infeasible sample to an operating point that is feasible to the N-1 security-constrained DC-OPF:

$$\min_{P_G^c, P_L, P_W, \theta^c} \|\mathbf{x} - \mathbf{x}^{\text{infeas}}\|_{\infty} \quad (21)$$

$$\text{s.t. (15) - (20)} \quad (22)$$

$$\mathbf{x} = [P_G^{\text{wosb}, 0}, P_W, P_L] \quad (23)$$

As this optimization problem is convex, we can provably identify the closest feasible sample \mathbf{x} to $\mathbf{x}^{\text{infeas}}$. Note that we solve this optimization problem to compute the results denoted with 'SC-DC-OPF' in Table II and Table IV.

C. N-1 Security-Constrained AC-OPF

In Section V-D and Section V-E, we create datasets of operating points classified according to their feasibility to the N-1 security constrained AC-OPF. In Section V-D, we consider small-signal stability constraints as well. These will be discussed in the Appendix A-D. We can formulate the N-1 security constrained AC-OPF problem as follows:

$$\min_{\mathbf{z}^c} f(\mathbf{P}_g^0) \quad (24)$$

$$\text{s.t. } \mathbf{z}^c := \{V^c, \theta^c, P_G^c, Q_G^c, P_W, P_L, Q_L\} \quad \forall c \in \mathcal{C} \quad (25)$$

$$\mathbf{s}_{\text{balance}}(\mathbf{z}^c) = \mathbf{0} \quad \forall c \in \mathcal{C} \quad (26)$$

$$|\mathbf{s}_{\text{line}}(\mathbf{z}^c)| \leq \mathbf{s}_{\text{line}}^{\max, c} \quad \forall c \in \mathcal{C} \quad (27)$$

$$P_G^{\min} \leq P_G^c \leq P_G^{\max} \quad \forall c \in \mathcal{C} \quad (28)$$

$$Q_G^{\min} \leq Q_G^c \leq Q_G^{\max} \quad \forall c \in \mathcal{C} \quad (29)$$

$$P_W^{\min} \leq P_W \leq P_W^{\max} \quad (30)$$

$$P_L^{\min} \leq P_L \leq P_L^{\max} \quad (31)$$

$$Q_L^{\min} \leq Q_L \leq Q_L^{\max} \quad (32)$$

$$\mathbf{V}^{\min} \leq \mathbf{V}^c \leq \mathbf{V}^{\max} \quad \forall c \in \mathcal{C} \quad (33)$$

$$P_G^{\text{wosb}, 0} = P_G^{\text{wosb}, c}, V_G^c = V_G^c \quad \forall c \in \mathcal{C} \quad (34)$$

The vector \mathbf{z}^c collects all variables for the intact and outaged system states in (25). The non-linear AC power flow nodal balance $\mathbf{s}_{\text{balance}}$ in (26) has to hold for intact and outaged system states. The absolute apparent line flow $|\mathbf{s}_{\text{line}}|$ is constrained by an upper limit in (27). For the full mathematical formulation, for brevity, please refer to [40]. The upper and lower limits on the system variables are defined in the constraints (28) – (33). The constraint (34) enforces preventive control actions for N-1 security: Both the generator active power set-points

and voltage set-points remain fixed during an outage. Note that the vector \mathbf{V}_G refers to the voltage set-points of the generators. We do not fix the entry in \mathbf{P}_G corresponding to the slack bus, as the slack bus generator compensates the mismatch in the losses. The independent variables that characterize an operating point are $\mathbf{x} := [\mathbf{P}_G^{\text{wosb},0}, \mathbf{P}_W, \mathbf{P}_L, \mathbf{Q}_L, \mathbf{V}_G^0]$. To create the datasets, based on the operating point \mathbf{x} , we solve the AC power flow equations in (26) to determine the dependent variables for each contingency and the intact system state. Then, we check the satisfaction of the operational constraints of the N-1 SC-AC-OPF problem including active and reactive generator limits (28) and (29), apparent branch flow limits (27), and voltage limits (33). For an operating point to be classified as feasible to the the N-1 SC-AC-OPF problem, it must satisfy all these constraints for all considered contingencies.

D. Small-Signal Stability

For the IEEE 14 bus test case in Section V-D we evaluate the feasibility of operating points with respect to combined small-signal stability and N-1 security. For the small signal stability model, we rely on standard system models and parameters commonly used for small signal stability analysis defined in Refs. [43] and [44]. We use a sixth-order synchronous machine model with an Automatic Voltage Regulator (AVR) Type I with three states. A more detailed description of the system model and parameters can be found in the Appendix of Ref. [45]. To determine small-signal stability, we linearize the dynamic model of the system around the current operating point, and compute the eigenvalues λ of the resulting system matrix \mathbf{A} . If all eigenvalues have negative real parts (i.e. lie in the left-hand plane), the system is considered small-signal stable, otherwise unstable. This can be formalized as follows:

$$\mathbf{A}(\mathbf{z}^c)\boldsymbol{\nu}^c = \boldsymbol{\lambda}^c \boldsymbol{\nu}^c \quad \forall c \in \mathcal{C} \quad (35)$$

$$\boldsymbol{\lambda}^c \leq 0 \quad \forall c \in \mathcal{C} \quad (36)$$

The set of variables \mathbf{z}^c is defined in (25). The term $\boldsymbol{\nu}^c$ denotes the right-hand eigenvectors of the system matrix \mathbf{A} . As we consider both N-1 security and small-signal stability, we have to modify the small-signal stability model for each operating point and contingency. We use Mathematica to derive the small signal model symbolically, MATPOWER AC power flows to initialize the system matrix, and Matlab to compute its eigenvalues and damping ratio, and assess the small-signal stability for each operating point and contingency.

REFERENCES

- [1] L. Duchesne, E. Karangelos, and L. Wehenkel, "Recent developments in machine learning for energy systems reliability management," *Proceedings of the IEEE*, 2020.
- [2] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *International Conference on Learning Representations (ICLR 2015)*, 2015.
- [3] L. A. Wehenkel, *Automatic learning techniques in power systems*. Springer Science & Business Media, 2012.
- [4] M. Glavic, "(deep) reinforcement learning for electric power system control and related problems: A short review and perspectives," *Annual Reviews in Control*, 2019.
- [5] L. Duchesne, E. Karangelos, and L. Wehenkel, "Using machine learning to enable probabilistic reliability assessment in operation planning," in *2018 Power Systems Computation Conference*. IEEE, 2018.
- [6] R. Dobbe, O. Sondermeijer, D. Fridovich-Keil, D. Arnold, D. Callaway, and C. Tomlin, "Towards distributed energy services: Decentralizing optimal power flow with machine learning," *IEEE Transactions on Smart Grid*, pp. 1–1, 2019.
- [7] S. Karagiannopoulos, P. Aristidou, and G. Hug, "Data-driven local control design for active distribution grids using off-line optimal power flow and machine learning techniques," *IEEE Transactions on Smart Grid*, vol. 10, no. 6, pp. 6461–6471, Nov 2019.
- [8] G. Dalal, E. Gilboa, S. Mannor, and L. Wehenkel, "Chance-constrained outage scheduling using a machine learning proxy," *IEEE Transactions on Power Systems*, 2019.
- [9] S. R. R. S. Kumar, and A. T. Mathew, "Online static security assessment module using artificial neural networks," *IEEE Transactions on Power Systems*, vol. 28, no. 4, pp. 4328–4335, Nov 2013.
- [10] V. J. Gutierrez-Martinez, C. A. Canizares, C. R. Fuente-Esquivel, A. Pizano-Martinez, and X. Gu, "Neural-network security-boundary constrained optimal power flow," *IEEE Transactions on Power Systems*, vol. 26, no. 1, pp. 63–72, Feb 2011.
- [11] F. Li and Y. Du, "From alphago to power system ai: What engineers can learn from solving the most complex board game," *IEEE Power and Energy Magazine*, vol. 16, no. 2, pp. 76–84, March 2018.
- [12] B. Donnot, I. Guyon, M. Schoenauer, P. Panciatichi, and A. Marot, "Introducing machine learning for power system operation support," *X Bulk Power Systems Dynamics and Control Symposium*, 2017.
- [13] B. Donnot, I. Guyon, M. Schoenauer, A. Marot, and P. Panciatichi, "Fast power system security analysis with guided dropout," *26th European Symposium on Artificial Neural Networks*, 2018.
- [14] M. Sun, I. Konstantelos, and G. Strbac, "A deep learning-based feature extraction framework for system security assessment," *IEEE Transactions on Smart Grid*, 2018.
- [15] J.-M. H. Arteaga, F. Hancharou, F. Thams, and S. Chatzivasileiadis, "Deep learning for power system security assessment," in *13th IEEE PowerTech 2019*. IEEE, 2019.
- [16] Y. Du, F. F. Li, J. Li, and T. Zheng, "Achieving 100x acceleration for n-1 contingency screening with uncertain scenarios using deep convolutional neural network," *IEEE Transactions on Power Systems*, 2019.
- [17] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, "The limitations of deep learning in adversarial settings," in *2016 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 2016, pp. 372–387.
- [18] Y. Chen, Y. Tan, and D. Deka, "Is machine learning in power systems vulnerable?" in *2018 IEEE SmartGridComm*. IEEE, 2018, pp. 1–6.
- [19] V. Tjeng, K. Y. Xiao, and R. Tedrake, "Evaluating robustness of neural networks with mixed integer programming," in *International Conference on Learning Representations (ICLR 2019)*, 2019.
- [20] C. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006.
- [21] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, pp. 436–444, 2015.
- [22] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, 2011, pp. 315–323.
- [23] H. Zhang, T.-W. Weng, P.-Y. Chen, C.-J. Hsieh, and L. Daniel, "Efficient neural network robustness certification with general activation functions," in *Advances in Neural Information Processing Systems 31*, 2018, pp. 4939–4948.
- [24] D. M. Kline and V. L. Berardi, "Revisiting squared-error and cross-entropy functions for training neural network classifiers," *Neural Computing & Applications*, vol. 14, no. 4, pp. 310–318, 2005.
- [25] M. Abadi et al., "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015, software available from tensorflow.org. [Online]. Available: <https://www.tensorflow.org/>
- [26] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, 2010, pp. 249–256.
- [27] M. Wu et al., "A game-based approximate verification of deep neural networks with provable guarantees," *Theoretical Computer Science*, 2019.
- [28] F. Croce and M. Hein, "Provable robustness against all adversarial l_p -perturbations for $p \geq 1$," *arXiv preprint arXiv:1905.11213*, 2019.
- [29] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," in *International Conference on Learning Representations (ICLR 2018)*, 2018.

- [30] F. Thams, A. Venzke, R. Eriksson, and S. Chatzivasileiadis, "Efficient database generation for data-driven security assessment of power systems," *IEEE Transactions on Power Systems*, 2019.
- [31] M. Zhu and S. Gupta, "To prune, or not to prune: exploring the efficacy of pruning for model compression," *International Conference on Learning Representations (ICLR 2018)*, 2018.
- [32] K. Dvijotham, R. Stanforth, S. Gowal, T. Mann, and P. Kohli, "A dual approach to scalable verification of deep networks," in *Conference on Uncertainty in Artificial Intelligence (UAI-18)*. Corvallis: AUAI Press, 2018, pp. 162–171.
- [33] K. Y. Xiao, V. Tjeng, N. M. M. Shafiullah, and A. Madry, "Training for faster adversarial robustness verification via inducing reLU stability," in *International Conference on Learning Representations (ICLR 2019)*, 2019.
- [34] R. D. Zimmerman, C. E. Murillo-Sánchez, and R. J. Thomas, "Matpower: Steady-state operations, planning, and analysis tools for power systems research and education," *IEEE Transactions on Power Systems*, vol. 26, no. 1, pp. 12–19, 2010.
- [35] M. D. McKay, R. J. Beckman, and W. J. Conover, "Comparison of three methods for selecting values of input variables in the analysis of output from a computer code," *Technometrics*, vol. 21, no. 2, pp. 239–245, 1979.
- [36] IEEE PES Task Force on Benchmarks for Validation of Emerging Power System Algorithms, "The power grid library for benchmarking AC optimal power flow algorithms," *arXiv:1908.02788*, Aug. 2019. [Online]. Available: <https://github.com/power-grid-lib/pglib-opf>
- [37] A. Venzke, D. K. Molzahn, and S. Chatzivasileiadis, "Efficient creation of datasets for data-driven power system applications," *XXI Power Systems Computation Conference (PSCC 2020)*, 2020.
- [38] D. Shchetinin, T. T. De Rubira, and G. Hug, "Efficient bound tightening techniques for convex relaxations of AC optimal power flow," *IEEE Transactions on Power Systems*, vol. 34, no. 5, pp. 3848–3857, 2019.
- [39] K. Sundar *et al.*, "Optimization-based bound tightening using a strengthened QC-relaxation of the optimal power flow problem," *arXiv:1809.04565*, 2018.
- [40] M. B. Cain, R. P. O'Neill, and A. Castillo, "History of optimal power flow and formulations," *Federal Energy Regulatory Commission*, vol. 1, pp. 1–36, 2012.
- [41] B. Stott, J. Jardim, and O. Alsac, "Dc power flow revisited," *IEEE Transactions on Power Systems*, vol. 24, no. 3, pp. 1290–1300, 2009.
- [42] F. Capitanescu, J. M. Ramos, P. Panciatici, D. Kirschen, A. M. Marcolini, L. Platbrood, and L. Wehenkel, "State-of-the-art, challenges, and future trends in security constrained optimal power flow," *Electric Power Systems Research*, vol. 81, no. 8, pp. 1731–1741, 2011.
- [43] P. W. Sauer and M. A. Pai, *Power system dynamics and stability*. Prentice hall Upper Saddle River, NJ, 1998, vol. 101.
- [44] F. Milano, *Power system modelling and scripting*. Springer Science & Business Media, 2010.
- [45] F. Thams, L. Halilbasic, P. Pinson, S. Chatzivasileiadis, and R. Eriksson, "Data-driven security-constrained opf," in *X Bulk Power Systems Dynamics and Control Symposium*, 2017.



Andreas Venzke (S'16) received the M.Sc. degree in Energy Science and Technology from ETH Zurich, Zurich, Switzerland in 2017. He is currently working towards the Ph.D. degree at the Department of Electrical Engineering, Technical University of Denmark (DTU), Kongens Lyngby, Denmark. His research interests include power system operation under uncertainty, convex relaxations of optimal power flow and machine learning applications for power systems.



Spyros Chatzivasileiadis (S'04, M'14, SM'18) is an Associate Professor at the Technical University of Denmark (DTU). Before that he was a post-doctoral researcher at the Massachusetts Institute of Technology (MIT), USA and at Lawrence Berkeley National Laboratory, USA. Spyros holds a PhD from ETH Zurich, Switzerland (2013) and a Diploma in Electrical and Computer Engineering from the National Technical University of Athens (NTUA), Greece (2007). In March 2016, he joined the Center for Electric Power and Energy at DTU. He is currently working on power system optimization and control of AC and HVDC grids, and machine learning applications for power systems.