

Adversarial Attacks on Time Series

Fazle Karim, *Graduate Student Member, IEEE*, Somshubra Majumdar[✉],
and Houshang Darabi[✉], *Senior Member, IEEE*

Abstract—Time series classification models have been garnering significant importance in the research community. However, not much research has been done on generating adversarial samples for these models. These adversarial samples can become a security concern. In this paper, we propose utilizing an adversarial transformation network (ATN) on a distilled model to attack various time series classification models. The proposed attack on the classification model utilizes a distilled model as a surrogate that mimics the behavior of the attacked classical time series classification models. Our proposed methodology is applied onto 1-nearest neighbor dynamic time warping (1-NN DTW) and a fully convolutional network (FCN), all of which are trained on 42 University of California Riverside (UCR) datasets. In this paper, we show both models were susceptible to attacks on all 42 datasets. When compared to Fast Gradient Sign Method, the proposed attack generates a larger fraction of successful adversarial black-box attacks. A simple defense mechanism is successfully devised to reduce the fraction of successful adversarial samples. Finally, we recommend future researchers that develop time series classification models to incorporating adversarial data samples into their training data sets to improve resilience on adversarial samples.

Index Terms—Time series classification, adversarial machine learning, perturbation methods, deep learning

1 INTRODUCTION

OVER the past decade, machine learning and deep learning have been powering several aspects of society [1]. Machine learning and deep learning are being used in some areas such as web searches [2], recommendation systems [3], and wearables [4]. With the advent of smart sensors, advancements in data collection and storage at vast scales, ease of data analytics and predictive modeling, time series data being collected from various sensors can be analyzed to determine regular patterns that are interpretable and exploitable. Classifying these time series data has been an area of interest by several researchers [5], [6], [7], [8]. Time series classification models are being used in health care, where ECG data are used to detect patients with severe cognitive defects, in audio, where words are categorized into different phenomes, and in gesture recognition, where motion data is used to categorize actions being made. Sensor data for resource and safety-critical applications such as manufacturing plants, industrial engineering, and chemical compound synthesis, when augmented by on-device analytics would allow automated response to avert significant issues in normal operation [9]. A successful time series classification model is able to capture and generalize the pattern of time series signals such that it is able to classify unseen data. Similarly, computer vision classification models exploit the spatial structure intrinsic to images obtained in the real world. However, computer vision models

have been shown to make incorrect predictions on the seemingly correct input, which is termed as an adversarial attack. Utilizing a variety of adversarial attacks, complex models are tricked to incorrectly predict the wrong class label. This is a serious security issue in neural networks widely used for vision-based tasks where adding slight perturbations or carefully crafted noise on an input image can mislead the image classification algorithm to make highly confident, yet wildly inaccurate predictions [10], [11]. This has been a growing concern in the Computer Vision field, where Deep Neural Networks (DNN) have been shown to be particularly susceptible to attacks [12], [13]. While DNNs are state-of-the-art models for a variety of classification tasks in several fields, including time series classification [14], [15], [16], these vulnerabilities harmfully impact real-world applicability in domains where secure and dependable predictions are of paramount importance [10], [17]. Compounding the severity of this issue, recent work by Papernot *et al.* has shown that adversarial attacks on a particular computer vision classifier can easily be transferred into other similar classifiers [18]. Only recently has the focus of attacks been shifted to Time Series classification models based on deep neural networks and classical models [12].

Several adversarial sample crafting techniques have been proposed to trick various image classification models that rely on DNN (state-of-the-art models for computer vision). Most of these techniques target the gradient information from the DNNs, which make them susceptible to these attacks [19], [20], [21]. Currently, new research is being done in generating natural language adversarial samples [22], [23]. This is extremely difficult due to the semantic-preserving perturbations. Even though time series classification models and natural language processing (NLP) models use similar deep learning modules (1D CNN, LSTM) the domains are completely different. Further, NLP models work better than time series classification models on discrete input data

- F. Karim and H. Darabi are with the Mechanical and Industrial Engineering, University of Illinois at Chicago, Chicago, IL 60607 USA.
E-mail: {karim1, hdarabi}@uic.edu.

- S. Majumdar is with the Computer Science, University of Illinois at Chicago, Chicago, IL 60607 USA. E-mail: smajum6@uic.edu.

Manuscript received 25 Feb. 2019; revised 27 Feb. 2020; accepted 29 Mar. 2020.
Date of publication 10 Apr. 2020; date of current version 2 Sept. 2021.
(Corresponding author: Houshang Darabi.)
Recommended for acceptance by D. Lowd.
Digital Object Identifier no. 10.1109/TPAMI.2020.2986319

because they model different patterns and structures. While adversarial attacks on NLP models which accept discrete tokens as input have been well studied, there has been relatively little study on adversarial attacks on time series models, which accept real values time series sequence as input. Generating adversarial samples for time series classification model has not been studied as much, in spite of the potentially large security risk they may possess. One major security concern exists in voice recognition tasks that convert speech-to-text. Carlini and Wagner [24] show how speech-to-text classifiers can be attacked. In addition, they provide various audio clips where a speech-to-text classifier, DeepSpeech, is not able to correctly detect the speech. Other security concerns can exist in healthcare devices that use time series classification algorithms, where it can be tricked into misdiagnosing patients that can affect the diagnosis of their ailment. Time series classification algorithms used to detect and monitor seismic activity can be manipulated to create fear and hysteria in our society. Wearables that use time series data to classify activity of the wearer can be fooled into convincing the users they are doing other actions. Most of the current state-of-the-art time series classification algorithms are classical approaches, such as 1 Nearest Neighbor - Dynamic Time Warping (1-NN DTW) [25], Kernel SVMs [26], or sophisticated methods such as Weasel [27], COTE [28], and Fast-Shapelet [29]. However, DNNs are fast becoming excellent time series classifiers due to their simplicity and effectiveness. The traditional time series classification models are harder to attack as it can be considered a black-box model with a non-differentiable internal computation. As such, no gradient information can be exploited. However, DNN models are more susceptible to white-box attacks as their gradient information can easily be exploited. A white-box attack is where the adversary is "given access to all elements of the training procedure" [21] - which includes the training dataset, training algorithm, the parameters and weights of the model, and the model architecture itself [21]. In other words, during a white-box attack, the attacker has full knowledge about the time series classification model (parameters, hyper-parameters, architecture, etc). Conversely, a black-box attack only has access to the target model's training procedure and model architecture [21]. During a black-box attack, the attacker has almost no knowledge of the time series classification model. In some cases, the only knowledge known when performing black-box attacks is the length of the input time series data. In this paper, we propose a black-box and a white-box attack that can attack both classical and deep learning time series classification state-of-the-art models.

In this work, we propose a proxy attack strategy on a target classifier via a student model, trained using standard model distillation techniques to mimic the behavior of the target classical time series classification models. The "student" network is the neural network distilled from another time series classification model, called the "teacher" model, that learns to approximate the output of the teacher model. Once the student model has been trained, our proposed adversarial transformation network (ATN) can then be trained to attack this student model. We apply our methodologies onto 1-NN DTW, Fully Connected Network and Fully Convolutional Network (FCN) that are trained on 42 University of California

Riverside (UCR) datasets [9]. When compared to a Baseline adversarial attack (Fast Gradient Sign Method), our proposed attack is able to generate a larger fraction of successful adversarial black-box attacks. Further, a simple defense mechanism is devised to reduce the fraction of successful adversarial samples on various time series classification attacks. Finally, we recommend future researchers that develop time series classification models to consider model robustness as an evaluative metric and incorporate adversarial data samples into their training data sets in order to further improve resilience to adversarial attacks.

The remainder of this paper is structured as follows: Section 2 provides a brief background on a couple time series classification models and background information on a few adversarial crafting techniques used on computer vision problems. Section 3 details our proposed methodologies and Section 4 presents and explains the results of our proposed methodologies on a couple of time series classification models. Section 5 concludes the paper and proposes future work.

2 BACKGROUND & RELATED WORKS

2.1 Time Series Classification Models

2.1.1 1-NN Dynamic Time Warping

The equations and definitions below are obtained from Kate *et al.* [30] and Xi *et al.* [8]. Dynamic Time Warping is a measure of similarity between 2 time series, Q and C , which is detected by finding their best alignment. Time series Q and C are defined as

$$Q = q_1, q_2, q_3, \dots, q_i, \dots, q_n \quad (1)$$

$$C = c_1, c_2, c_3, \dots, c_i, \dots, c_m. \quad (2)$$

To align both the time series data, the distance between each timestep of Q and C is calculated, $(q_i - c_j)^2$, to generate a n -by- m matrix. In other words, the i th and j th of the matrix is the q_i and c_j . The optimal alignment between Q and C is considered the warping path, W , such that $W = w_1, w_2, w_3, \dots, w_k, \dots, w_K$. The warping path is computed such that,

- 1) $w_1 = (1, 1)$,
- 2) $w_k = (n, m)_k$,
- 3) given $w_k = (a, b)$ then $w_{k-1} = (a', b')$ where $0 \leq a - a' \leq 1$ and $0 \leq b - b' \leq 1$.

The optimal alignment is the warping path that minimizes the total distance between the aligning points

$$DTW(Q, C) = \arg \min_{W=w_1, w_2, \dots, w_K} \sqrt{\sum_{k=1, w_k=(i,j)}^k (q_i - c_j)^2}. \quad (3)$$

2.1.2 Fully Convolutional Network

The Fully Convolutional Network (FCN) is one of the earliest deep learning time series classifier [31]. It contains 3 convolutional layers, with convolution kernels of size 8, 5 and 3 respectively, and emitting 128, 256 and 128 filters respectively. Each convolution layer is followed by a batch normalization [32] layer that is applied with a ReLU activation layer. A global average pooling layer is employed after the

last ReLU activation layer. Finally, softmax is applied to determine the class probability vector.

2.2 Adversarial Transformation Network

Several methods have been proposed to generate adversarial samples that attack deep neural networks that are trained for computer vision tasks. Most of these methods use the gradient with respect to the image pixels of these neural networks. Baluja and Fischer [33] propose Adversarial Transformation Networks (ATNs) to efficiently generate an adversarial sample that attacks various networks by training a feed-forward neural network in a self-supervised method. Given the original input sample, ATNs modify the classifier outputs slightly to match the adversarial target. ATN works similarly to the generator model in the Generative Adversarial Training framework.

According to Baluja and Fischer *et al.* [33], an ATN can be parametrized as a neural network $g_f(x) : x \rightarrow \hat{x}$, where f is the target model (either a classical model or another neural network) which outputs either a probability distribution across class labels or a sparse class label, and $\hat{x} \sim x$, but $\text{argmax } f(x) \neq \text{argmax } f(\hat{x})$. To find g_f , we minimize the following loss function :

$$L = \beta * L_x(g_f(\mathbf{x}_i), \mathbf{x}_i) + L_y(f(g_f(\mathbf{x}_i)), f(\mathbf{x}_i)), \quad (4)$$

where L_x is a loss function on the input space (e.g., L_2 loss function), L_y is the specially constructed loss function on the output space of f to avoid learning the identity function, \mathbf{x}_i is the i th sample in the dataset and β is the weighing term between the two loss functions.

It is necessary to carefully select the loss function L_y on the output space to successfully avoid learning the identity function. Baluja and Fischer *et al.* [33] define the loss function L_y as $L_y(\mathbf{y}', \mathbf{y}) = L_2(\mathbf{y}', r(\mathbf{y}, t))$, where $\mathbf{y} = f(x)$, $\mathbf{y}' = f(g_f(x))$, t is the index of the target class such that $t \in [1 \dots C]$, where C is the number of classes and $r(\cdot)$ is a reranking function that modifies \mathbf{y} such that $y_k < y_t, \forall k \neq t$. This reranking function $r(\mathbf{y}, t)$ can either be the simple one hot encoding function $\text{onehot}(t)$ or be formulated to take advantage of the already present \mathbf{y} to encourage better reconstruction. We therefore utilize the reranking function proposed by Baluja and Fischer *et al.* [33], which can be formulated as

$$r_\alpha(\mathbf{y}, t) = \text{norm} \left(\begin{cases} \alpha * \text{maxy}, & \text{if } k = t \\ y_k, & \text{otherwise} \end{cases} \right)_{k \in \mathbf{y}}, \quad (5)$$

where $\alpha > 1$ is an additional hyperparameter which defines how much larger y_t should be than the current max classification and norm is a normalizing function that rescales its input to be a valid probability distribution

2.3 Transferability Property

Papernot *et al.* [18] propose a black-box attack by training a local substitute network, s , to replicate or approximate the target DNN model, f . The local substitute model is trained using synthetically generated samples and the output of these samples are labels from f . Subsequently, s is used to generate adversarial samples that it misclassifies. Generating adversarial samples for s is much easier, as its full knowledge/parameters are available, making it susceptible to various

attacks. The key criteria to successfully generate adversarial samples of f is the transferability property, where adversarial samples that misclassify s will also misclassify f .

2.4 Knowledge Distillation

Knowledge distillation, first proposed by Bucila *et al.* [34], is a model compression technique where a small model, s , is trained to mimic a pre-trained model, f . This process is also known as the model distillation training where the teacher is f and the student is s . The knowledge that is distilled from the teacher model to the student model is done by minimizing a loss function, where the objective of the student model is to imitate the distribution of the class probabilities of the teacher model. Hinton *et al.* [35] note that there are several instances where the probability distribution is skewed such that the correct class probability would have a probability close to 1 and the remaining classes would have a probability closer to 0. Hence, Hinton *et al.* [35] recommend computing the probabilities q_i from the pre-normalized logits z_i , such that

$$q_i = \sigma(z; T) = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}, \quad (6)$$

where T is a temperature factor normally set to 1. Higher values of T produce softer probability distributions over classes. The loss that is minimized is the model distillation loss, further explained in Section 3.3.

3 METHODOLOGY

3.1 Gradient Adversarial Transformation Network

In this work, we employ distinct methodologies for white-box and black-box attacks, in order to adhere to a strictly realistic set of limitations in black-box attacks. For both methodologies, we incorporate Adversarial Transformation Networks (ATN) [33] as a generative neural network that accepts an input time series sample x and transforms it to an adversarial sample \hat{x} .

An Adversarial Transformation Network can be formulated as a neural network $g_f(x) : x \rightarrow \hat{x}$, where f is the model that will be attacked. We further augment the information available to the ATN with the gradient of the input sample x with respect to the softmax scaled logits of the target class predicted by the attacked classifier. We can therefore formally define a Gradient Adversarial Transformation Network (GATN) as a neural network parametrized as $g_f(x, \tilde{x}) : [x; \tilde{x}] \rightarrow \hat{x}$, where

$$\tilde{x} = \frac{\partial f_t}{\partial x}, \quad (7)$$

such that $x \in \mathbb{R}^T$ is an input time series of maximum length T , f_t represents the probability of the input time series being classified as the target class, t , and $[;]$ represents the concatenation operation of two vectors. GATN computes the gradient of the input with respect to the output as the objective of an adversarial network is to learn the perturbations necessary to alter the input in order to affect the classification outcome. Hence, with the availability of the input gradient \tilde{x} , the Gradient Adversarial Transformation Network can better construct adversarial samples that can affect the targeted model while reducing the overall perturbation added

to the sample. Therefore we utilize the GATN model for all of our attacks.

A significant issue with the above formulation of the GATN is the non-differentiability of classical models. Distance-based models such as 1-NN Dynamic Time Warping do not have the notion of gradients during either training or evaluation. Therefore, we cannot directly compute the gradient of the input (\hat{x}) with respect to the 1-NN DTW model f . We discuss the solution to this issue in Section 3.3, by building a student neural network s which approximates the predictions of the non-differentiable classical classifier f .

3.2 Black-Box and White-Box Restrictions

While this formulation of the GATN is sufficient for white-box attacks where we have access to the attacked model f or the student model s , this assumption is unrealistic in the case of black-box attacks. For a black-box, we are not permitted access to either the internal model (a neural network or a classical model) or to the dataset that the model was trained on. Furthermore, for black-box attacks, we impose a restriction on the predicted labels, such that we utilize only the class label predicted, and not the probability distribution produced after softmax scaling (for neural networks), or the scaled probabilistic approximations of classical model predictions.

To further restrict ourselves to realistic attack vectors, we stratify the available dataset D , which will be used to train the GATN, into two halves, such that we train the GATN on one subset, D_{eval} , and are able to perform evaluations on both this train set and the wholly unseen test set, D_{test} . Note that this available dataset D is not the dataset on which the attacked model f was trained on. As such, we never utilize the train set available to the attacked classifier to either train or evaluate the GATN model. In order to satisfy these constraints on available data, we define our available dataset D as the test set of the UCR Archive [9]. As the test set was never used to train any attacked model f , it is sufficient to utilize it as an unseen dataset. We then split the test dataset into two class-balanced halves, D_{eval} and D_{test} . Another convenience is the availability of test set labels, which can be harnessed as a strict check when evaluating adversarial generators.

When we evaluate under the constraints of black-boxes, we further limit ourselves to unlabeled train sets, where we assume the available dataset is unlabeled, and thereby utilize only the predicted label from the attacked classifier f to label the dataset prior to attacks. We state this as an important restriction, considering that it is far more difficult to freely obtain or create datasets for time series than for images which are easily understood and interpreted. For time series, significant expertise may be required to distinguish one sample amongst multiple classes, whereas natural images can be coarsely labeled with relative ease without sophisticated equipment or expertise.

3.3 Training Methodology

A chief consideration during training of ATN or GATN is the loss formulation on the prediction space (L_y) is heavily influenced by the reranking function $r(\cdot)$ chosen. If we opt for the one hot encoding of the target class, we lose the ability to maintain class ordering and the ability to adjust the ranking weight (α) to obtain adversaries with less distortion.

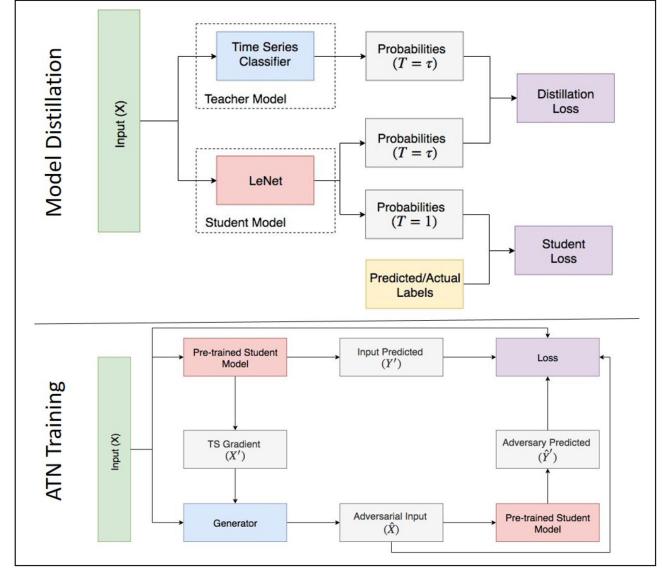


Fig. 1. The top diagram shows the methodology of training the model distillation used in the white-box and black-box attacks. The bottom diagram is the methodology utilized to attack a time series classifier.

However, to utilize the appropriate reranking function, we must have access to the class probability distribution, which is unavailable to black-box attacks. It may not even be possible to compute for certain classical models such as 1-NN DTW which uses distance-based computations to determine the nearest neighbor.

To overcome this limitation, we employ knowledge distillation as a mechanism to train a student neural network s , which is trained to replicate the predictions of the attacked model f . As such, we are required to compute the predictions of the attacked model on the dataset we possess just one time, which can be either class labels or probability distribution over all classes. We then utilize these labels as the ground truth labels that the student s is trained to imitate. In case the predictions are class labels, we utilize one hot encoding scheme to compute the cross entropy loss, otherwise, we try to imitate the probability distribution directly. It is to be noted that the student model shares the training dataset D_{eval} with the GATN model.

As suggested by Hinton *et al.* [35], we describe the training scheme of the student as shown in Fig. 1. We scale the logits of the student s and teacher f (iff the teacher provides probabilities and it is a white-box attack) by a temperature scaling parameter τ , which is kept constant at 10 for all experiments. When training the student model, we minimize the loss function defined as

$$L_{transfer} = \gamma * L_{distillation} + (1 - \gamma) * L_{student} \quad (8)$$

$$L_{distillation} = \mathcal{H}(\sigma(z_f; T = \tau), \sigma(z_s; T = \tau)) \quad (9)$$

$$L_{student} = \mathcal{H}(y, \sigma(z_s; T = 1)), \quad (10)$$

where \mathcal{H} is the standard cross entropy loss function, z_s and z_f are the un-normalized logits of the student (s) and teacher (f) models respectively, $\sigma(\cdot)$ is the scaled-softmax operation described in Equation (6), y is the ground truth labels, and γ is a gating parameter between the two losses and is used to maintain a balance between how much the

student s imitates the teacher f versus how much it learns from the hard label loss. When training a student as a white-box attack, we set γ to be 0.5, allowing the equal weight to both losses, whereas for a black-box attack, we set γ to be 1. Therefore for black-box attacks, we force the student s to only mimic the teacher f to the limit of its capacity. In setting this restriction, we limit the amount of information that may be made available to the GATN.

Once we have a student model s which is capable of simulating the predictions of the attacked model f , we then train the GATN using this student model. Fig. 1 shows the methodology of training such a model. Since the GATN requires not just the original sample x but also the gradient of that sample \tilde{x} with respect to the predictions for the targeted class, we require two forward passes from the student model. The first forward pass is simply to obtain the gradient of the input \tilde{x} , as well as the predicted probability distribution of the student y . The adversarial sample crafted (\hat{X}) is then used in a second forward pass to compute the predicted probability distribution of the student with respect to the adversarial sample, \hat{y}' . We minimize the weighted loss measure L defined in Section 2.2 in order to train the GATN model.

3.4 Evaluation Methodology

Due to the different restrictions imposed between available information depending on whether the attack is a white-box or black-box attack, we train the GATN on one of two models. We assert that we train the GATN by attacking the target neural network f directly only when we perform a white-box attack on a neural network. In all other cases, whether the attack is a white-box or black-box attack, and whether the attacked model is a neural network or a classical model, we select the student model s as the model which is attacked to train the GATN, and then use the GATN's predictions (\hat{x}) to check if the teacher model f is also attacked when provided the predicted adversarial input (\hat{x}) as a sample.

During evaluation of the trained GATN, we compute the number of adversaries of the attacked model f that have been obtained on the training set D_{eval} . During the evaluation, we can measure any metric under two circumstances. Provided a labeled dataset which was split, we can perform a two-fold verification of whether an adversary was found or not. First, we check that the ground truth label matches the predicted label of the classifier when provided with an unmodified input ($y = y'$ when input x is provided to f), and then check whether this predicted label is different from the predicted label when provided with the adversarial input ($y \neq y'$ when input \hat{x} is provided to f). This ensures that we do not count an incorrect prediction from a random classifier as an attack.

Another circumstance is that we do not have any labeled samples prior to splitting the dataset. This training set is an unseen set for the attacked model f , therefore we consider that the dataset is unlabeled, and assume that the label predicted by the base classifier is the ground truth ($y = y'$ by default, when sample x is provided to f). This is done prior to any attack by the GATN and is computed just once. We then define an adversarial sample as a sample \hat{x} whose predicted class label is different than the predicted ground truth label ($y \neq y'$, when sample \hat{x} is provided to f). A drawback of this approach is that it is overly optimistic and

rewards sensitive classifiers that misclassify due to very minor alterations.

In order to adhere to an unbiased evaluation, we chose the first option, and utilize the provided labels that we know from the test set to properly evaluate the adversarial inputs. In doing so, we acknowledge the necessity of a labeled test set, but as shown above, it is not strictly necessary to follow this approach.

4 EXPERIMENTS AND RESULTS

All methodologies were tested on 42 benchmark datasets for time series classification found in the UCR repository. The 42 datasets selected were all from the types "Sensor", "ECG", "EOG", and "Hemodynamics", where an adversarial attack is a potential security concern. The 42 datasets contain data with application varying from classifying chlorine concentration to earthquakes. The ECG datasets contain a few time series classification problems that require classification of humans with heart conditions or Myocardial Infarctions. A couple of the EOG and Hemodynamics datasets require the classification of Japanese Katakana strokes and heart air pressure of pigs respectively. Further detailed information on these datasets can be found online [9]. These 42 datasets are all the datasets in all the domains in the repository that possess a security concern. With the exception of images and motion, we believe the remaining domains do not possess a serious security concern realistically. Adversarial attacks in the domain of images and motion is well studied and is the reason why the proposed time series classification adversarial methodologies are not used upon it.

We evaluate based on two criterion, the mean squared error between the training dataset and the generated samples (lower is better) and the fraction of successful adversaries (higher is better). For all experiments, we keep α , the reranking weight, set to 1.5, the target class set to 1, and perform a grid search over 5 possible values of β , the reconstruction weight term, such that $\beta = 10^{-b}; b \in \{1, 2, 3, 4, 5\}$. In addition, all student models are trained only using D_{eval} . The codes and weights of all models are available at https://github.com/houshd/TS_Adv.

4.1 Experiments

We select both neural networks as well as traditional models as the attacked model f . For the attacked neural network, we utilize a Fully Convolutional Network, whereas for the base traditional model, 1NN-Dynamic Time Warping Classifier is utilized.

To maintain the strictest definition of the black and white-box attacks, we utilize only the discrete class label of the attacked model for black-box attacks and utilize the probability distribution predicted by the classifier for white-box attacks. The only exception where a student-teacher network is not used is when performing a white-box attack on a FCN time series model, as the gradient information of a neural network can be directly exploited by an Adversarial Transformation Network (ATN). The performance of the adversarial model is evaluated on the original time series classification teacher model.

For every student model we train, we utilize the LeNet-5 architecture [36]. The student models are trained only using

D_{eval} . We define a LeNet-5 time series classifier as a classical Convolutional Neural network following the structure : Conv (6 filters, 5×5 , valid padding) — Max Pooling — Conv (16 filters, 5×5 , valid padding) — Max Pooling — Fully Connected (120 units, relu) — Fully Connected (84 units, relu) — Fully Connected (number of classes, softmax).

The fully convolutional network is based on the FCN model proposed by Wang *et al.* [31]. It is comprised of 3 blocks, each comprised of a sequence of Convolution layer - Batch Normalization—ReLU activations. All convolutional kernels are initialized using the uniform he initialization proposed by He *et al.* [37]. We utilize [128, 256, 128] filters and kernel sizes of [8, 5, 3] to be consistent.

A strong deterministic baseline model to classify time series is 1-NN DTW with 100 percent warping window. Due to its reliance on a distance matrix as a means of its classification, it cannot easily be used to compute an equivalent soft probabilistic representation. Since white-box attacks have access to the probability distribution predicted for each sample, we utilize this distance matrix in the computation of an equivalent soft probabilistic representation. The equivalent representation is such that if we compute the top class (class with highest probability score) on this representation, we get the exact same result as selecting the 1-nearest neighbor on the actual distance matrix.

To compute this soft probabilistic representation, consider a distance matrix V computed using a distance measure such as DTW between all possible pairs of samples between the two datasets being compared.

Algorithm 1. Equivalent Probabilistic Representation of the Distance Matrix for 1-Nearest Neighbor Classification

```

1 Algorithm: Soft-1NN ( $V, y$ )
2 Data:  $V$  is a distance matrix of shape  $[N_{test}, N_{train}]$  and  $y$  is
   the train set label vector of length  $N_{train}$ 
3 Result: Softmax normalized predictions  $p$  of shape  $[N_{test}, C]$ 
   and the discrete label vector  $q$  of length  $N_{test}$ 
4 begin
5    $V \leftarrow (-V)$ 
6    $uniqueClasses = Unique(y)$  // class labels
7    $V_c = []$ 
8   for  $c_i$  in  $uniqueClasses$  do
9      $v_c = V_{(y=c_i)} // [N_{test}, N_{train}(y = c_i)]$ 
10     $v_{c\_max} = max(v_c) // [N_{test}]$ 
11     $V_c.append(v_{c\_max})$ 
12  end
13   $V' = concatenate(V_c) // [N_{test}, number of classes]$ 
14   $p = softmax(V') // [N_{test}, number of classes]$ 
15   $q = argmax(p) // [N_{test}]$ 
16  return ( $p, q$ )
17 end
```

Algorithm 1 is an intermediate normalization algorithm which accepts a distance matrix V and the class labels of the training set y as inputs and computes an equivalent probabilistic representation that can directly be utilized to compute the 1-nearest neighbor. The *Soft-1NN* algorithm selects all samples that belong to a class c_i , where $i \in \{1, \dots, C\}$ as v_c , computes the maximum over all train samples for that class and appends the vector v_{c_max} to the list V_c . The concatenation of all of these lists of vectors in V_c then

represents the matrix V' , on which we then apply the *softmax* function, as shown in Equation (6) with T set to 1, to represent this matrix V' as a probabilistic equivalent of the original distance matrix V .

An implicit restriction placed on Algorithm 1 is that the representation is equivalent only when computing the 1-nearest neighbor. It cannot be used to represent the K -nearest neighbors and therefore cannot be used for K -nearest neighbor classification. However, in time series classification, the general consensus is on the use of 1-nearest neighbor classifiers and its variants to classify time series [6], [7], [25], [28], [38]. While the above algorithm has currently been applied to convert the 1INN-DTW distance matrix, it can also be applied to normalize any distance matrix utilized for 1-NN classification algorithms.

4.2 Results

Figs. 2 and 3 depict the results from white-box attacks on 1-NN DTW and FCN that is applied on 42 UCR datasets. Further, Figs. 4 and 5 represent the results from black-box attacks on 1-NN DTW and FCN classifiers that are trained on the same 42 UCR datasets. The detailed results can be found in Appendix A, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2020.2986319>. The proposed methodology is successfully in capturing adversaries on all datasets. An example of an adversarial attack on the dataset "FordB" is shown in Fig. 6.

4.2.1 Fraction of Successful Adversaries

The fraction of successful adversaries (the number of successful adversaries divided by the total number of correctly classified samples by the original classifier) and amount of perturbation per sample in each dataset can increase or decrease depending on the hyper-parameters that are tested on. For example, the dataset "Trace" has 0 adversaries for most of the attacks (black-box attack on 1-NN DTW, white-box attack on 1-NN DTW, black-box attack on FCN) when the Target Class is set to 1. However, if the target class is changed to 2, the number of adversaries generated increases to 9,3,1,37 for a black-box attack on 1-NN DTW, white-box attack on 1-NN DTW, black-box attack on FCN and white-box attack on FCN, respectively. These numbers can be higher if the hyper-parameters are changed. In addition, due to the loss function of the ATN, the target class has a significant impact on the adversary being generated. It is easier to generate adversaries for time series classes that are similar to each other.

A Wilcoxon signed-rank test is utilized to compare the fraction of successful adversaries generated by white-box and black-box attacks on FCN and 1-NN classifiers that are trained on the 42 datasets, summarized in Table 1. Our results indicate that the FCN classifier is more susceptible to a white-box attack compared to a white-box attack on 1-NN DTW. It is to be noted that the white-box attack on the FCN classifier generates significantly more fraction of successful adversaries than its counterparts. This is because the white-box attack is directly on the FCN model and not on a student model that approximates the classifier behavior. We observe that the fraction of successful adversarial samples

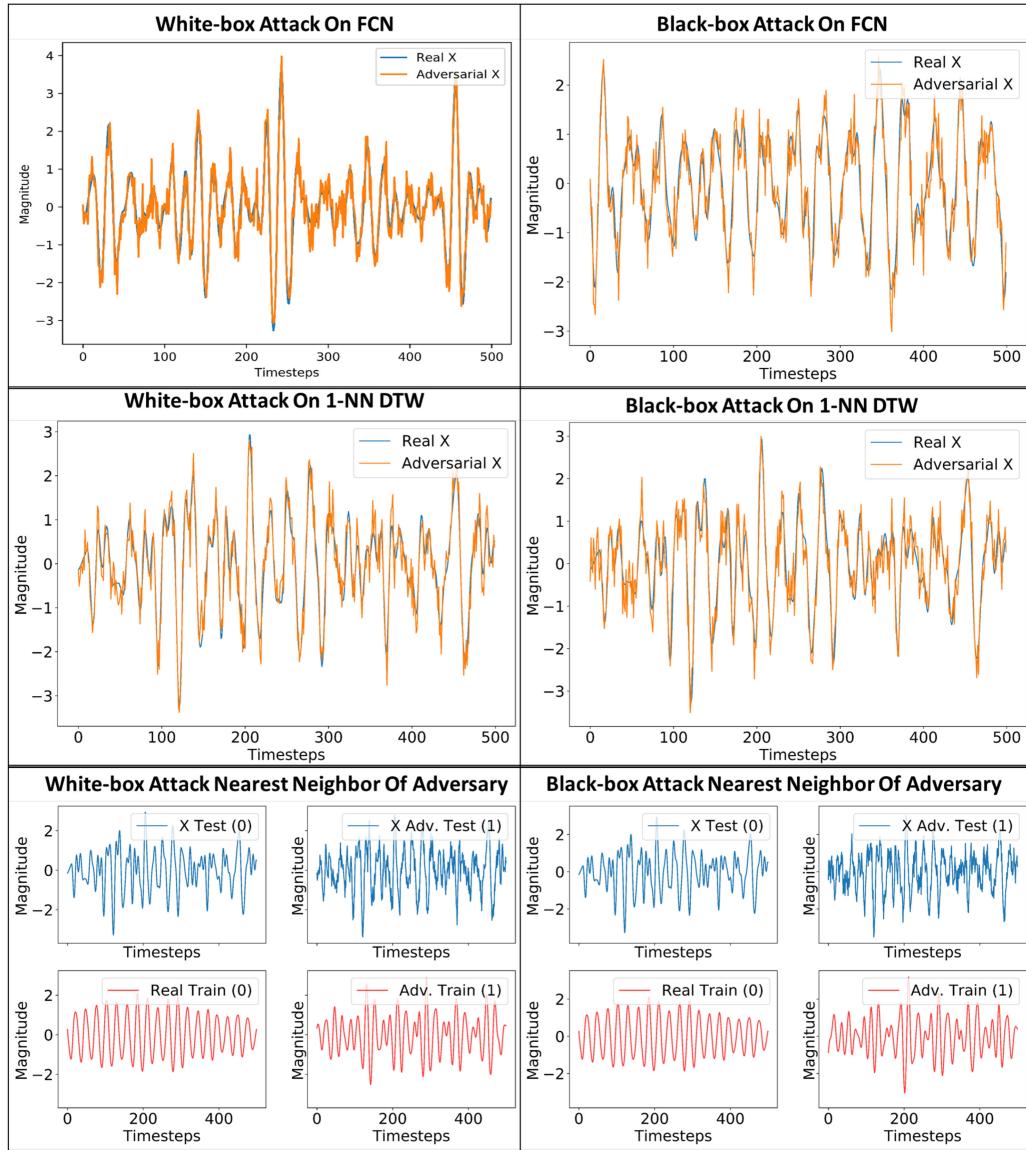


Fig. 6. A sample black-box and white-box attack on an FCN and 1-NN DTW classifier that is trained on the dataset “FordB”. The last row of the figure depicts the nearest neighbor of the original and adversarial time series.

4.2.2 Comparison to a Baseline Model

The proposed architectures are compared to a baseline adversarial attack, Fast Gradient Sign Method (FGSM) [39]. FGSM requires the gradient of each model. The black-box attacks and attacks on classical time series classification models is difficult when using it as it does not have access to its gradient. However, this is mitigated when applying it on the student network. We apply a Wilcoxon signed-rank test to compare the fraction of successful adversaries generated by FGSM to the fraction of successful adversaries generated by the respective black-box or white-box attack on FCN or 1-

NN DTW using GATN. This is summarized in Table 2. The black-box attack on FCN and 1-NN DTW using GATN generates significantly more adversaries than when applying FGSM. However, FGSM generates significantly more adversaries than GATN when performing a white-box attack on FCN. A white-box attack on FCN using GATN performs statistically the same as a white-box attack on 1-NN DTW classifiers using FGSM. This indicates FGSM works better when the actual gradients of the classifier are accessible. GATN works better on black-box attacks.

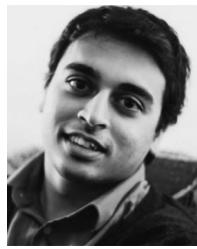
4.2.3 Mean Squared Error (MSE)

We use MSE as a metric to depict how much perturbation is required for each dataset. Since, each dataset is normalized with a mean of 0 and a variance of 1, a lower MSE is better as each sample requires a lesser amount of perturbation. Hence, the optimal scenario would be when the MSE is closer to 0. The average MSE of adversarial samples after

TABLE 1
Wilcoxon Signed-Rank Test Comparing the Fraction of Successful Adversarial Between the Different Attacks

	White-box 1-NN DTW	Black-box FCN	White-box FCN
Black-box 1-NN DTW	2.278E-01	9.850E-02	5.949E-08
White-box 1-NN DTW		1.345E-01	2.731E-07
Black-box FCN			3..198E-03

- [18] N. Papernot, P. McDaniel, and I. Goodfellow, "Transferability in machine learning: From phenomena to black-box attacks using adversarial samples," 2016, *arXiv:1605.07277*.
- [19] N. Akhtar and A. Mian, "Threat of adversarial attacks on deep learning in computer vision: A survey," *IEEE Access*, vol. 6, pp. 14410–14430, 2018.
- [20] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *stat*, vol. 1050, p. 9, 2017.
- [21] F. Tramèr, A. Kurakin, N. Papernot, I. Goodfellow, D. Boneh, and P. McDaniel, "Ensemble adversarial training: Attacks and defenses," *stat*, vol. 1050, p. 30, 2017.
- [22] M. Alzantot, Y. Sharma, A. Elgohary, B.-J. Ho, M. Srivastava, and K.-W. Chang, "Generating natural language adversarial examples," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2018.
- [23] R. Jia and P. Liang, "Adversarial examples for evaluating reading comprehension systems," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2017, pp. 2021–2031.
- [24] N. Carlini and D. Wagner, "Audio adversarial examples: Targeted attacks on speech-to-text," in *IEEE Secur. Privacy Workshops*, pp. 1–7, 2018.
- [25] E. Keogh and C. A. Ratanamahatana, "Exact indexing of dynamic time warping," *Knowl. Inf. Syst.*, vol. 7, no. 3, pp. 358–386, 2005.
- [26] A. Kampouraki, G. Manis, and C. Nikou, "Heartbeat time series classification with support vector machines," *IEEE Trans. Inf. Technol. Biomed.*, vol. 13, no. 4, pp. 512–518, Jul. 2009.
- [27] P. Schäfer and U. Leser, "Fast and accurate time series classification with weasel," in *Proc. ACM Conf. Inf. Knowl. Manage.*, 2017, pp. 637–646.
- [28] A. Bagnall, J. Lines, J. Hills, and A. Bostrom, "Time-series classification with COTE: The collective of transformation-based ensembles," *IEEE Trans. Knowl. Data Eng.*, vol. 27, no. 9, pp. 2522–2535, Sep. 2015.
- [29] T. Rakthanmanon and E. Keogh, "Fast shapelets: A scalable algorithm for discovering time series shapelets," in *Proc. SIAM Int. Conf. Data Mining*, 2013, pp. 668–676.
- [30] R. J. Kate, "Using dynamic time warping distances as features for improved time series classification," *Data Mining Knowl. Discov.*, vol. 30, no. 2, pp. 283–312, 2016.
- [31] Z. Wang, W. Yan, and T. Oates, "Time series classification from scratch with deep neural networks: A strong baseline," in *Proc. Int. Joint Conf. Neural Netw.*, 2017, pp. 1578–1585.
- [32] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," 2015, *arXiv:1502.03167*.
- [33] S. Baluja and I. Fischer, "Adversarial transformation networks: Learning to generate adversarial examples," 2017, *arXiv:1703.09387*.
- [34] C. Buciluă, R. Caruana, and A. Niculescu-Mizil, "Model compression," in *Proc. 12th ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining*, 2006, pp. 535–541.
- [35] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *stat*, vol. 1050, p. 9, 2015.
- [36] Y. LeCun *et al.*, "LeNet-5, convolutional neural networks," 2015, Art. no. 20. [Online]. Available: <http://yann.lecun.com/exdb/lenet>
- [37] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2015, pp. 1026–1034.
- [38] P. Schäfer, "The boss is concerned with time series classification in the presence of noise," *Data Mining Knowl. Discov.*, vol. 29, no. 6, pp. 1505–1530, 2015.
- [39] X. Yuan, P. He, Q. Zhu, and X. Li, "Adversarial examples: Attacks and defenses for deep learning," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 30, no. 9, pp. 2805–2824, Sep. 2019.



Fazle Karim (Graduate Student Member, IEEE) received the BSc degree in industrial engineering from the University of Illinois at Urbana-Champaign, Champaign, Illinois, in 2012, the MSc degree in industrial engineering from the University of Illinois at Chicago, Chicago, Illinois, in 2016. He is currently working toward the PhD degree from the Mechanical and Industrial Engineering Department, University of Illinois at Chicago, Chicago, Illinois. He is also the lead data scientist with Prominent Laboratory, the university's foremost research facility in process mining. His research interests include education data mining, health care data mining, and time series analysis.



Somshubra Majumdar received the BS degree in computer engineering from the University of Mumbai, Mumbai, Maharashtra, India, in 2016, and the MS degree in computer science from the University of Illinois at Chicago, Chicago, Illinois, in 2018. He is currently an aspiring artificial intelligence researcher. His research interests include the domain of image classification and segmentation using convolutional neural networks, time series classification using recurrent neural networks, speech recognition, and machine learning.



Houshang Darabi (Senior Member, IEEE) received the PhD degree in industrial and systems engineering from Rutgers University, New Brunswick, New Jersey, in 2000. He is currently an associate professor with the Department of Mechanical and Industrial Engineering, University of Illinois at Chicago (UIC), and also an associate professor with the Department of Computer Science, UIC. He has been a contributing author of two books in the areas of scalable enterprise systems and reconfigurable discrete event systems. His research has been supported by several federal and private agencies, such as the National Science Foundation, the National Institute of Standard and Technology, the Department of Energy, and Motorola. He has extensively published on various automation and project management subjects, including wireless sensory networks for location sensing, planning and management of projects with tasks requiring multi-mode resources, and workflow modeling and management. He has published in different prestigious journals and conference proceedings, such as the *IEEE Transaction on Robotics and Automation*, *IEEE Transactions on Automation Science and Engineering*, and *IEEE Transactions on Systems, Man, and Cybernetics*, and the *Information Sciences*. His current research interests include the application of data mining, process mining, and optimization in design and analysis of manufacturing, business, project management, and workflow management systems.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/csdl.