

Real-Time Adversarial Attacks, by Yuan Gong et.al.

Presented by Ramesh Kr Sah

March 18, 2020

Table of Contents

1 Introduction

2 Problem Definition

3 Solution

4 Conclusion

Real-Time Adversarial Attacks

- ❶ All existing adversarial attack methods work with complete input sample to find adversarial perturbations.
- ❷ These attack approaches are not applicable to situations where the target model takes streaming input.
- ❸ In this paper, the author proposes a real-time adversarial attack scheme for target model with streaming inputs, such as speech recognition system.
- ❹ The proposed attack continuously uses observed data to approximate an optimal adversarial perturbation for future time points using RNN.

Real-Time Adversarial Attack Scheme

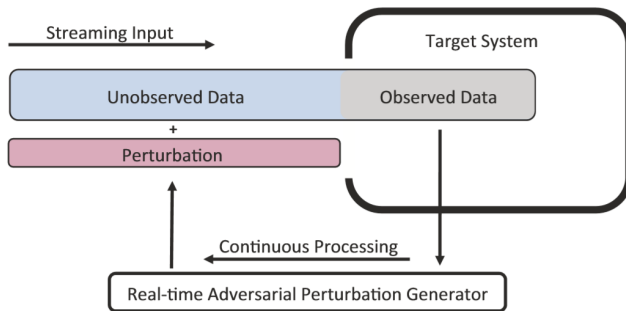


Figure 1: An illustration of the real-time adversarial attack scheme. The target system takes streaming input; only past data points can be observed and adversarial perturbation can only be added to future data points. The adversarial perturbation generator continuously uses observed data to approximate an optimal adversarial perturbation for future data points.

Problem Formalization

For time-series input $x = \{x_1, \dots, x_n\} \in \mathbb{R}^{m \times n}$, where each point $x_i \in \mathbb{R}^m$, and classifier $f : x \rightarrow y$, where $y \in \{1, \dots, k\}$ output labels.

The attacker designs a real-time adversarial perturbation generator $g(\cdot)$ to generate the adversarial perturbation r_t such that,

$$r_t = \begin{cases} g(\{x_1, x_2, \dots, x_{t-d-1}\}) & d+1 < t < n \\ 0 & \text{else} \end{cases} \quad (1)$$

And with a metric $m(\cdot)$ to measure perceptibility of adversarial perturbations, it solves the optimization problem:

$$\begin{aligned} &\text{minimize} && m(r = \{r_1, r_2, \dots, r_n\}) \\ &\text{s.t} && f(x + r) \neq f(x) \end{aligned} \quad (2)$$

Equation 1 makes sure that adversarial perturbation is crafted only based on the observed part of the data sample and can only be applied to the unobserved part.

And equation 2 constraints the perturbation to be as imperceptible as possible on the premise that the attack succeeds.

Approach

The adversarial perturbation generator $g(.)$ is an agent which tries to solve the *partially observable decision process* problem. The problem is described using tuple $\langle O, S, A, T, R \rangle$, where

- ① **Observation O:** $o_t = \{x_1, x_2, \dots, x_t\}$
- ② **State S:** unobservable hidden state.
- ③ **Action A:** $a_t = r_{t+d+1}$ i.e., adding the perturbation to original sample at time $t + d + 1$.
- ④ **Transition T :** unknown
- ⑤ **Reward R:** $I_{f(x+r) \neq f(x)} - m(r)$

The goal of the RL is to learn an optimal policy $\pi_g : a_t = g(o_t)$ that maximizes the expectation of the reward.

Why RL Perturbation Generator?

- 1 Since the input sample obeys some fixed distribution, the perturbation generator will be able to forecast future perturbations for unobserved data.
- 2 Also, dependencies among the data points of the data sample, will provide learning opportunities for the generator to help forecast perturbations for future data points based on the observed data points.
- 3 And, with RL we can utilize the state-of-the-art non-real-time adversarial generation algorithms to generate trajectories of observation-action pairs, to learn the optimal trajectory for the perturbation generator.

Use Imitation Learning to Train the Agent

In imitation learning an RL agent learns an optimal policy π_g by imitating the behavior of an expert. The expert is a state-of-the-art non-real-time adversarial example crafting technique as the expert.

- ① Use the expert to generate sample-perturbation pairs $\langle (\mathbf{x}^1, \mathbf{r}^1), (\mathbf{x}^2, \mathbf{r}^2), \dots \rangle$ as decision trajectories by feeding different original sample \mathbf{x}^i and collecting the corresponding adversarial perturbations \mathbf{r}^i .
- ② Convert the \mathbf{x} and \mathbf{r} into \mathbf{o} and \mathbf{a} , and build the dataset D .
- ③ By treating \mathbf{o} as the input feature and \mathbf{a} as the output label, learn $\pi_g : \mathbf{a}_t = g(\mathbf{o}_t)$ in a supervised manner.

Once we form the dataset D , consisting of observation-action pairs from the expert's decision trajectory, the real-time perturbation generator g , a deep neural network, is trained using D .

Algorithms - Getting the Trajectories

Algorithm 1 Real-time Adversarial Attack

Require:

Original dataset $\mathcal{X} = \{x^i\}$ where each $x^i = \{x_t^i\}$

Non-real-time adversarial example generator (expert) g_e

Phase 1: Generate Expert Demonstrations

Input: Original sample set \mathcal{X}

Output: Expert decision trajectory set \mathcal{D}

- 1: initialize \mathcal{D} as an empty set
- 2: **for** each $x^i \in \mathcal{X}$ **do**
- 3: $r^i = \{r_t^i\} = g_e(x^i)$
- 4: initialize trajectory τ_i as an empty set
- 5: **for** each time point t of x^i **do**
- 6: $o_t^i = \{x_1^i, x_2^i, \dots, x_t^i\}$
- 7: $a_t^i = r_{t+d+1}^i$
- 8: add (o_t^i, a_t^i) to τ_i
- 9: **end for**
- 10: add τ_i to \mathcal{D}
- 11: **end for**
- 12: **return** \mathcal{D}

Algorithms - Training the Agent and Evaluation

Phase 2: Train Realtime Adversarial Example Generator

Input: Expert decision trajectory set \mathcal{D}

Output: Real-time adversarial example generator g_r

- 13: initialize g_r as a recurrent network with parameter θ
- 14: **for** each trajectory $\tau_i \in \mathcal{D}$ **do**
 - ▷ maintain RNN states for each t to expedite computing
- 15: **for** each time point t **do**
 - 16: calculate the predicted action $\hat{a}_t^i = g_r(o_t^i)$
 - 17: calculate the loss l between the predicted action \hat{a}_t^i and the expert's action a_t^i
 - 18: update θ to minimize the loss l
- 19: **end for**
- 20: **end for**
- 21: **return** g_r

Phase 3: Conduct Real-time Adversarial Attack

Input: Streaming observations $\mathbf{o} = \{o_1, o_2, \dots, o_t, \dots\}$

- 22: **at** each time point t **do**
 - 23: $a_t = g_r(o_t)$
 - 24: execute action a_t
-

Model Details

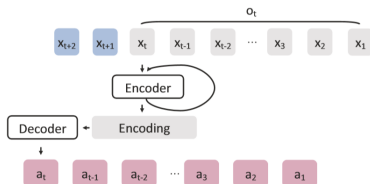


Figure 3: Illustration of the training process. Note that the output action only depends on the current observation o_t .

- 1 The generator has two components: encoder and decoder.
- 2 The encoder is a recurrent neural network that maps a variable length input into a fixed dimensional encoding.
- 3 The decoder then makes the decision of the action, and calculates the error between the predicted action and the ground truth action to update g using back-prop.

Model Architecture

Layer Name	Output Dimension
Input	$(t, 1)$
Framing	$(\lceil t/160 \rceil, 160)$
Conv1 / Pooling	$(\lceil t/160 \rceil, 80, 16)$
Conv2 / Pooling	$(\lceil t/160 \rceil, 40, 32)$
Conv3 / Pooling	$(\lceil t/160 \rceil, 20, 48)$
Conv4 / Pooling	$(\lceil t/160 \rceil, 10, 64)$
Flatten	$(\lceil t/160 \rceil, 640)$
LSTM * 3	(256)
Dense 1	(256)
Dense 2	(128)
Output	(5)

Table 1: The network details and the output dimension of each layer.

Attacking a Voice Command Recognition System

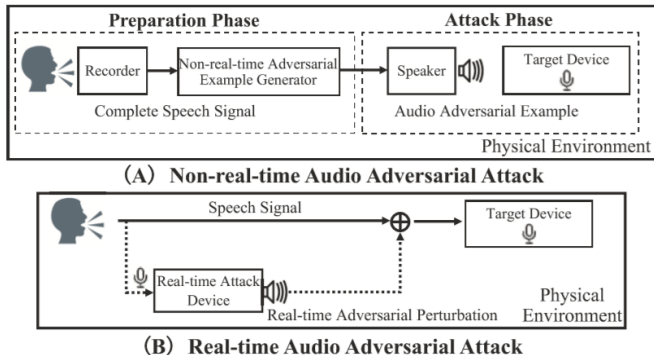


Figure 4: Illustration of the non-real-time (upper figure) and real-time (lower figure) audio adversarial attack.

Results

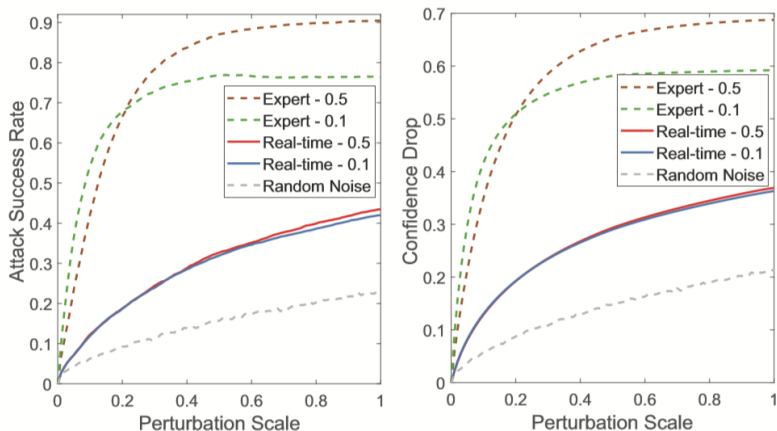


Figure 7: The attack successful rate and the confidence score dropped by the attack with different perturbation scale.

- ① Real-Time Adversarial Attack Method, but the performance is not at par to offline attack methods.
- ② Performance may possibly improved using sophisticated and efficient algorithms for perturbation generator.
- ③ Extend this work by incorporating the knowledge from mistakes made by the agent to facilitate learning.

Thank You!

Link to the paper: <https://www.ijcai.org/Proceedings/2019/0649.pdf>