

Data-Driven Attack Compensate Control of Neural Networks

1st Qidong Liu

*School of Automation Engineering
University of Electronic Science and Technology of China
Chengdu, China*

Abstract—In this paper, we propose a novel data-driven approach to compensate for the effects of cyber attacks on neural network-based control systems. As neural networks become increasingly integral to critical control applications, these systems face heightened vulnerability to adversarial attacks. Our approach utilizes historical data to theoretically analyze and compensate for deviations in control performance caused by such attacks. The method integrates attack detection with a compensation mechanism designed to adjust the control input in real-time, aiming to mitigate the impact of the attack. Through rigorous theoretical analysis, we demonstrate the potential of this approach to enhance system stability and performance in the presence of cyber threats.

Index Terms—Data-driven control, Neural networks, Cyber attack compensation, Control systems, Attack detection.

I. INTRODUCTION

The adoption of neural networks (NN) in control systems has seen a rapid increase due to their ability to approximate complex nonlinear functions and their flexibility in learning from data. However, this widespread adoption has also made these systems susceptible to various forms of cyber attacks. Adversarial attacks on neural networks, such as perturbation-based attacks and data poisoning, can significantly degrade system performance or even lead to catastrophic failures in safety-critical applications.

Given the high stakes involved, it is crucial to develop robust control strategies that can detect and compensate for the effects of such attacks in real time. Traditional control methods often rely on pre-defined models and assumptions, which may not be valid under attack scenarios. In contrast, data-driven approaches offer a promising alternative by utilizing real-time data to adapt the control strategy dynamically.

A. Relevance of Data-Driven Control

Data-driven control methods have gained prominence due to their ability to handle complex, nonlinear, and uncertain systems without requiring an explicit mathematical model. These methods rely on historical and real-time data to infer the system's behavior and adjust the control inputs accordingly. In the context of neural network-based control systems, a data-driven approach can be particularly effective in compensating for the unforeseen effects of cyber attacks, as it can learn and adapt to the changes in system dynamics caused by the attack.

Recent research has explored various data-driven techniques for control and estimation in the presence of uncertainties [1],

[2]. However, the application of these techniques specifically to compensate for cyber attacks on neural networks remains relatively unexplored. This paper seeks to fill this gap by proposing a data-driven control framework that integrates attack detection with compensation mechanisms.

B. Problem Statement and Objectives

Consider a control system where the control input $u(t)$ is generated by a neural network controller $\mathcal{N}(x(t); \theta)$, where $x(t)$ is the state vector, and θ represents the neural network parameters. The system dynamics can be expressed as:

$$\dot{x}(t) = f(x(t)) + g(x(t))u(t),$$

where $f(\cdot)$ and $g(\cdot)$ are nonlinear functions representing the system dynamics. In the presence of an attack, the control input may be compromised, leading to a modified input $u_a(t) = u(t) + \Delta u(t)$, where $\Delta u(t)$ represents the perturbation caused by the attack.

The objective of this paper is to design a data-driven compensation mechanism that can detect the presence of the attack and adjust the control input $u(t)$ to mitigate the impact of $\Delta u(t)$. Specifically, we aim to achieve the following:

1. **Attack Detection:** Develop a mechanism that utilizes historical and real-time data to detect deviations in control performance indicative of an attack.
2. **Compensation Control:** Design a compensation strategy that adjusts the control input based on the detected attack, ensuring that the system remains stable and performs within acceptable bounds.

II. ATTACK DETECTION MECHANISM

To detect the presence of an attack, we propose a data-driven anomaly detection algorithm that monitors the residual $r(t)$ between the expected system output $\hat{y}(t)$ generated by the neural network model and the actual system output $y(t)$. The residual is defined as:

$$r(t) = y(t) - \hat{y}(t).$$

Under normal operating conditions, $r(t)$ is expected to be small and bounded. However, when an attack occurs, $r(t)$ is likely to deviate significantly from its normal range. By analyzing the statistical properties of $r(t)$, we can identify the onset of an attack.

A threshold ϵ is set based on the historical distribution of $r(t)$, and an attack is detected if $|r(t)| > \epsilon$ for a specified duration. This detection mechanism is integrated with the compensation control strategy to ensure that corrective actions are taken promptly.

III. COMPENSATION CONTROL STRATEGY

Once an attack is detected, the control input $u(t)$ is adjusted to compensate for the perturbation $\Delta u(t)$. The compensation is achieved by introducing an adjustment term $\Delta_c u(t)$ to the original control input:

$$u_c(t) = u(t) + \Delta_c u(t),$$

where $\Delta_c u(t)$ is determined based on the estimated attack impact. We propose a data-driven approach to estimate $\Delta u(t)$ by using a neural network-based estimator that learns the relationship between the residual $r(t)$ and the required compensation $\Delta_c u(t)$.

The estimator is trained using historical attack data, and during operation, it continuously updates its parameters using real-time data to improve its accuracy. The overall compensation control law can be expressed as:

$$u_c(t) = u(t) - \alpha \hat{\Delta} u(t),$$

where $\hat{\Delta} u(t)$ is the estimated perturbation, and α is a gain parameter that controls the level of compensation.

A. Stability Analysis

To ensure the stability of the compensated control system, we analyze the closed-loop system dynamics with the compensation mechanism in place. The compensated system can be described by:

$$\dot{x}(t) = f(x(t)) + g(x(t))(u(t) - \alpha \hat{\Delta} u(t)).$$

Using a Lyapunov function $V(x(t)) = x(t)^T P x(t)$, where P is a positive definite matrix, we can derive conditions under which the system remains stable despite the presence of an attack. The derivative of $V(x(t))$ along the trajectories of the system is given by:

$$\dot{V}(x(t)) = x(t)^T \left(\frac{\partial f(x(t))}{\partial x} + \frac{\partial g(x(t))}{\partial x} (u(t) - \alpha \hat{\Delta} u(t)) \right) x(t).$$

By appropriately choosing the compensation gain α and ensuring that $\hat{\Delta} u(t)$ accurately reflects the attack impact, we can guarantee that $\dot{V}(x(t))$ is negative definite, thus ensuring stability.

REFERENCES

- [1] S. M. Dibaji, M. Pirani, D. B. Flamholz, A. M. Anaswamy, K. H. Johansson, and A. Chakraborty, "A Systems and Control Perspective of CPS Security," *Annual Reviews in Control*, vol. 47, pp. 394-411, 2019.
- [2] Y. Wang, A. Shafiee, and M. Althoff, "Verification of Neural Network Controlled Systems with Simplex Architecture," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 2, pp. 722-732, 2021.
- [3] N. Carlini, D. Wagner, "Towards Evaluating the Robustness of Neural Networks," *IEEE Symposium on Security and Privacy*, pp. 39-57, 2017.
- [4] X. He, C. Zhang, and S. Ren, "Robustness of Deep Neural Networks to Adversarial Attacks," *IEEE Access*, vol. 6, pp. 15983-15991, 2018.