

# Optimal Tracking Control of Linear Discrete-Time Systems Under Cyber Attacks<sup>\*</sup>

Hao Liu,<sup>\*</sup> Hui Qiu

<sup>\*</sup> School of Automation, Shenyang Aerospace University, Shenyang, 110136, China (e-mail: lh\_hit\_1985@163.com).

---

**Abstract:** In this paper, the optimal tracking control problem is solved based on the reinforcement learning for linear systems subject to multiple false-data-injection (FDI) attacks. An augmented system is established, which includes the original system and reference-trajectory generator system. The corresponding optimal control issue is formulated as a game problem between the system and malicious adversaries. A Q-learning algorithm is proposed to solve the game algebraic Riccati equation without requiring any knowledge about the dynamics of the augmented system. Finally, an example is provided to show that the system output can track the reference trajectory under cyber attacks.

*Keywords:* Linear systems, optimal tracking control; false-data-injection attacks; game theory; Q-learning.

---

## 1. INTRODUCTION

Optimal tracking control is a main research field in control theory, and the objective is to design an optimal controller such that the system output can track a reference trajectory in an optimal sense. This can be achieved by minimizing a predefined quadratic performance index. It should be noted that the communication is vulnerable to attacks launched by malicious adversaries when the controller and the plant are linked through wireless network. Thus, it is meaningful to consider the effect of cyber-attacks.

In Teixeira et al. (2015), cyber-secure networked control is modeled and analyzed, and the attack space is divided into three dimensions: the adversary's model knowledge, disclosure and disruption resources. Consequently, the corresponding attacks are mainly classified into three types, i.e., denial-of-service attacks Persis et al. (2016); Qin et al. (2018); Ding et al. (2017), false-data-injection (FDI) attacks Kung et al. (2017); Bai et al. (2017); Hu et al. (2018) and replay attacks Mo et al. (2009); Zhu et al. (2014); Chen et al. (2018). FDI attacks aim to replace the original data with false data injected by adversaries. Note that these injected data may deteriorate the system performance. Moreover, FDI attacks require to know the model knowledge and disruption resources.

The security problem for a networked control system was investigated in Hu et al. (2018), and a new necessary and sufficient condition for the insecurity was derived. A notion of  $\epsilon$ -stealthiness was given in Bai et al. (2017) to quantify the detectability of attacks in stochastic cyber-physical systems. Based on the entropy theory, the performance degradation in the presence of attacks was characterized.

---

<sup>\*</sup> This work is supported by the National Natural Science Foundation of China (61703286) and the Project from the Education Department of Liaoning Province (JYT19039).

It is noted that a Stackelberg game theory can be used to deal with FDI attacks if the two players make their decisions sequentially. In Li et al. (2018), the interactive decision-making process between the defender and the attacker were studied in a Stackelberg game framework, and the optimal strategies for both sides were solved by using linear programming approach.

It should be emphasized that, however, the design of controller is rarely considered to defend FDI attacks. For instance, the controller design problem of CPS was studied in Ye et al. (2019) to ensure the reliability and security when actuator faults in physical layers and attacks in cyber layers occur simultaneously. In Ding et al. (2017), the consensus control problem was investigated for a class of multiagent systems with lossy sensors and cyber-attacks, in which a dynamic output feedback controller was designed.

As a result, optimal tracking control for discrete-time systems under multiple FDI attacks are investigated in this work. It is assumed that the controller can acquire the information of the plant state and the reference signal. Moreover, the corresponding issue is formulated into a game problem, which can be solved by using reinforcement learning method.

*Notations:* The superscript  $T$  stands for matrix transposition;  $\mathbb{R}^n$  denotes the  $n$ -dimensional Euclidean space;  $I$  and  $0$  represent the identity matrix and the zero matrix, respectively; the notation  $P > 0$  means that  $P$  is real symmetric and positive definite;  $A \otimes B$  is the Kronecker product of  $A$  and  $B$ .

## 2. PROBLEM FORMULATION

Consider the following discrete-time linear systems

$$x_{k+1} = Ax_k + Bu_k \quad (1)$$

$$y_k = Cx_k \quad (2)$$

where  $x_k \in \mathbb{R}^n$  is the system state,  $u_k \in \mathbb{R}^m$  is the control input, and  $y_k \in \mathbb{R}^p$  is the system output. Assume that  $(A, B)$  and  $(A, C)$  are controllable and observable pairs, respectively.

In this paper, it is assumed that the reference input satisfies the following model

$$y_{k+1}^r = Ty_k^r \quad (3)$$

where  $y_k^r \in \mathbb{R}^p$  and  $T$  needs not to be Hurwitz. The aim of controller is to make the system output  $y_k$  track the reference input  $y_k^r$ . Moreover, we assume that the plant and the controller are linked through wireless network. Assume that the wireless network communication can be attacked by adversaries and all attackers have full knowledge of the system. The communication between the remote controller and the system can be eavesdropped and attacked, and the modified data  $u_k^a$  can be represented by

$$u_k^a = u_k + \sum_{j=1}^q \Gamma^j a_k^j \quad (4)$$

where  $q$  is the number of adversaries,  $a_k^j \in \mathbb{R}^m$ ,  $j \in \mathcal{Q} = \{1, 2, \dots, q\}$ , is the false-data injected by attacker  $j$  at time step  $k$ , and matrix  $\Gamma^j$  satisfies

$$\Gamma^j = \text{diag}\{\beta_1^j, \beta_2^j, \dots, \beta_m^j\} \quad (5)$$

where  $\beta_i^j \in \{0, 1\}$ ,  $i = 1, 2, \dots, m$ ,  $j \in \mathcal{Q}$ .  $\beta_i^j = 1$  means that  $i^{\text{th}}$  communication channel is attacked by  $j^{\text{th}}$  attacker; otherwise, the corresponding channel is not injected false data. It should be pointed out that  $\Gamma^j$  can be determined by using different method, such as game theory approach, assuming  $\beta_i^j$  satisfies bernoulli distribution, etc. In this paper, it is assumed that the matrix  $\Gamma^j$  is known in advance.

Considering modified control input (4), the equation (1) becomes

$$x_{k+1} = Ax_k + Bu_k^a = Ax_k + Bu_k + \sum_{j=1}^q B\Gamma^j a_k^j \quad (6)$$

Now, define the following tracking error

$$e_k = y_k - y_k^r = Cx_k - y_k^r \quad (7)$$

Based on (6) and the reference input dynamics (3), the following augmented system can be constructed

$$\begin{aligned} \bar{x}_{k+1} &= \begin{bmatrix} A & 0 \\ 0 & T \end{bmatrix} \bar{x}_k + \begin{bmatrix} B \\ 0 \end{bmatrix} u_k + \sum_{j=1}^q \begin{bmatrix} B\Gamma^j \\ 0 \end{bmatrix} a_k^j \\ &= \bar{A}\bar{x}_k + \bar{B}u_k + \sum_{j=1}^q \bar{\Gamma}^j a_k^j \end{aligned} \quad (8)$$

where  $\bar{x}_k = [x_k^T, (y_k^r)^T]^T$ . Then, the tracking error  $e_k$  can be given by

$$e_k = [C \quad -I]\bar{x}_k = \bar{C}\bar{x}_k \quad (9)$$

Note that the controller can be designed in the form of different types, such as state feedback, dynamic output feedback, etc. On the other hand, the false-data injected by adversaries can also be given in many different types, which means that an attacker can utilizes some different types of signals, such as the system states, the reference

input, and so on. In this paper, however, the control input  $u_k$  and the injected false data  $a_k^j$  are assumed to be linear functions of  $x_k$  and  $y_k^r$ , which can be represented as follows

$$u_k = \mu(x_k, y_k^r) = K_1 x_k + K_2 y_k^r = K \bar{x}_k \quad (10)$$

and

$$a_k^j = h^j(x_k, y_k^r) = L_1^j x_k + L_2^j y_k^r = L^j \bar{x}_k \quad (11)$$

where  $j \in \mathcal{Q}$ ,  $K$  and  $L^j$  are constant feedback gain matrices to be determined. The advantage of such design can be seen in the following analysis.

The objective of the defender, i.e., the controller, is to minimize the following reward function at time step  $k$

$$\mathcal{R}_d(x_k, y_k^r, u_k) = \sum_{i=k}^{\infty} \gamma^{i-k} (e_i^T Q_e e_i + u_i^T R u_i) \quad (12)$$

where  $Q_e > 0$ ,  $R > 0$  and  $\gamma \in (0, 1)$  is the discount factor. It is noted that  $\gamma = 1$  can be chosen if matrix  $T$  is Hurwitz. Therefore, the optimal control policy can be obtained by solving

$$u_k^* = \arg \min_{u_k} \mathcal{R}_d(x_k, y_k^r, u_k) \quad (13)$$

Similarly, the reward for the attacker  $j$  at time step  $k$  is

$$\mathcal{R}_a(x_k, y_k^r, a_k^j) = \sum_{i=k}^{\infty} \gamma^{i-k} (e_i^T Q_e e_i - \vartheta^j (a_i^j)^T (a_i^j)) \quad (14)$$

where  $\vartheta^j > 0$  is a pre-defined weighting parameter, and it is assumed that both players know the reward functions of their opponents. The optimal attack policy for  $j^{\text{th}}$  attacker can be given by

$$a_k^{j,*} = \arg \max_{a_k^j} \mathcal{R}_a(x_k, y_k^r, a_k^j), \quad j \in \mathcal{Q} \quad (15)$$

*Remark 1.* Equations (13) and (15) show that the defender aims to reduce the tracking error and consume less energy; while the purpose of  $j^{\text{th}}$  attacker is to increase the tracking error. Therefore, the objectives of both players are opposite.

### 3. MAIN RESULTS

Define the following functions

$$\begin{aligned} J(\bar{x}_k, u_k, a_k^1, \dots, a_k^q) &= J(x_k, y_k^r, u_k, a_k^1, \dots, a_k^q) \\ &= \sum_{i=k}^{\infty} \gamma^{i-k} (e_i^T Q_e e_i + u_i^T R u_i - \sum_{j=1}^q \vartheta^j (a_i^j)^T (a_i^j)) \end{aligned} \quad (16)$$

$$\bar{J} = \inf_{u_k} \sup_{a_k^1, \dots, a_k^q} J(\bar{x}_k, u_k, a_k^1, \dots, a_k^q) \quad (17)$$

$$\underline{J} = \sup_{a_k^1, \dots, a_k^q} \inf_{u_k} J(\bar{x}_k, u_k, a_k^1, \dots, a_k^q) \quad (18)$$

If  $J^* = \bar{J} = \underline{J}$ , then  $J^*$  is called the value of the corresponding game and  $(u_k^*, a_k^{1,*}, \dots, a_k^{q,*})$  is the Nash equilibrium at time step  $k$ , namely, the saddle-point solution.

Then, calculating (13) and (15) is equivalent to find the saddle point  $(u_k^*, a_k^{1,*}, \dots, a_k^{q,*})$  such that

$$\begin{aligned} J^*(x_k, y_k^r) &= \min_{u_k} \max_{a_k^1, \dots, a_k^q} J(x_k, y_k^r, u_k, a_k^1, \dots, a_k^q) \\ &= \max_{a_k^1, \dots, a_k^q} \min_{u_k} J(x_k, y_k^r, u_k, a_k^1, \dots, a_k^q) \end{aligned} \quad (19)$$

Define the following utility function at time step  $k$

$$\begin{aligned} r_k &= r_k(x_k, y_k^T, u_k, a_k^1, \dots, a_k^q) \\ &= \bar{x}_k^T Q_x \bar{x}_k + u_k^T R u_k - \sum_{j=1}^q \vartheta^j (a_k^j)^T (a_k^j) \end{aligned} \quad (20)$$

where  $Q_x = \bar{C}^T Q_e \bar{C}$ . Considering the state feedback control, then one can obtain from (16) that

$$\begin{aligned} J(\bar{x}_k) &= J(\bar{x}_k, \mu(\bar{x}_k), h^1(\bar{x}_k), \dots, h^q(\bar{x}_k)) \\ &= \sum_{i=k}^{\infty} \gamma^{i-k} r_k(\bar{x}_k, \mu(\bar{x}_k), h^1(\bar{x}_k), \dots, h^q(\bar{x}_k)) \end{aligned} \quad (21)$$

Through a simple calculation, (21) can be re-written as

$$\begin{aligned} J(\bar{x}_k) &= r_k(\bar{x}_k, \mu(\bar{x}_k), h^1(\bar{x}_k), \dots, h^q(\bar{x}_k)) + \gamma J(\bar{x}_{k+1}) \\ &= \bar{x}_k^T Q_x \bar{x}_k + u_k^T R u_k - \sum_{j=1}^q \vartheta^j (a_k^j)^T (a_k^j) + \gamma J(\bar{x}_{k+1}) \end{aligned} \quad (22)$$

which is the Bellman equation of optimal control. According to the Bellman optimality principle, the value function  $J(\bar{x}_k)$  is quadratic in the state  $\bar{x}_k$  at time step  $k$ , and can be represented as

$$J(\bar{x}_k) = \bar{x}_k^T P \bar{x}_k \quad (23)$$

where  $P \in \mathbb{R}^{(n+p) \times (n+p)}$  is positive definite and to be determined. Then, the Bellman equation (22) becomes

$$\begin{aligned} \bar{x}_k^T P \bar{x}_k &= \bar{x}_k^T Q_x \bar{x}_k + u_k^T R u_k - \sum_{j=1}^q \vartheta^j (a_k^j)^T (a_k^j) \\ &\quad + \gamma \bar{x}_{k+1}^T P \bar{x}_{k+1} \end{aligned} \quad (24)$$

For simplicity, we can define  $\vartheta = \text{diag}\{\vartheta^1, \vartheta^2, \dots, \vartheta^q\}$ ,  $\bar{\Gamma} = [\bar{\Gamma}^1 \bar{\Gamma}^2 \dots \bar{\Gamma}^q]$  and  $a_k = [(a_k^1)^T (a_k^2)^T \dots (a_k^q)^T]^T$ . Therefore, the optimal strategies satisfy

$$\begin{aligned} &\begin{bmatrix} R + \gamma \bar{B}^T P \bar{B} & \gamma \bar{B}^T P \bar{\Gamma} \\ \gamma \bar{\Gamma}^T P \bar{B} & \gamma \bar{\Gamma}^T P \bar{\Gamma} - \vartheta \otimes I \end{bmatrix} \begin{bmatrix} u_k \\ a_k \end{bmatrix} \\ &= - \begin{bmatrix} \gamma \bar{B}^T P \bar{A} \\ \gamma \bar{\Gamma}^T P \bar{A} \end{bmatrix} \bar{x}_k \end{aligned} \quad (25)$$

The following optimal policies for both players can be obtained by solving (25)

$$u_k^* = K(P) \bar{x}_k = (R + \gamma \bar{B}^T P \bar{B} - \Omega \bar{B})^{-1} \Omega \bar{A} \bar{x}_k \quad (26)$$

and

$$a_k^* = L(P) \bar{x}_k = -(\Theta \bar{\Gamma} - \vartheta \otimes I)^{-1} \Theta \bar{A} \bar{x}_k \quad (27)$$

where

$$\begin{aligned} \Omega &= \gamma^2 \bar{B}^T P \bar{\Gamma} (\gamma \bar{\Gamma}^T P \bar{\Gamma} - \vartheta \otimes I)^{-1} \bar{\Gamma}^T P \\ \Theta &= [(\Theta^1)^T (\Theta^2)^T \dots (\Theta^q)^T]^T \\ \Theta^j &= \gamma (\bar{\Gamma}^j)^T P [I - \gamma \bar{B} (R + \gamma \bar{B}^T P \bar{B})^{-1} \bar{B}^T P] \\ L(P) &= [(L^1(P))^T (L^2(P))^T \dots (L^q(P))^T]^T \end{aligned}$$

and  $P$  satisfies the following game Riccati equation

$$\begin{aligned} P &= Q_x + \gamma \bar{A}^T P \bar{A} - \gamma^2 [\bar{A}^T P \bar{B} \quad \bar{A}^T P \bar{\Gamma}] \\ &\quad \times \begin{bmatrix} R + \gamma \bar{B}^T P \bar{B} & \gamma \bar{B}^T P \bar{\Gamma} \\ \gamma \bar{\Gamma}^T P \bar{B} & \gamma \bar{\Gamma}^T P \bar{\Gamma} - \vartheta \otimes I \end{bmatrix}^{-1} \begin{bmatrix} \bar{B}^T P \bar{A} \\ \bar{\Gamma}^T P \bar{A} \end{bmatrix} \end{aligned} \quad (28)$$

In order to obtain the optimal strategies for both sides, the knowledge of system dynamics  $(A, B, \Gamma, T)$  is needed. In this paper, a Q-function is given for the game between the defender and attackers, which does not need to know

$(A, B, \Gamma, T)$ . Define the following Q-function Bellman equation

$$\begin{aligned} Q(\bar{x}_k, u_k, a_k^1, \dots, a_k^q) &= r_k(\bar{x}_k, u_k, a_k^1, \dots, a_k^q) + \gamma J(\bar{x}_{k+1}) \\ &= \bar{x}_k^T Q_x \bar{x}_k + u_k^T R u_k - \sum_{j=1}^q \vartheta^j (a_k^j)^T (a_k^j) + \gamma J(\bar{x}_{k+1}) \\ &= \bar{x}_k^T Q_x \bar{x}_k + u_k^T R u_k - a_k^T (\vartheta \otimes I) a_k \\ &\quad + \gamma Q(\bar{x}_{k+1}, u_{k+1}, a_{k+1}^1, \dots, a_{k+1}^q) \end{aligned} \quad (29)$$

which can be re-written as the following compact form

$$\begin{aligned} Q(\bar{x}_k, u_k, a_k^1, \dots, a_k^q) &= Q(\bar{x}_k, u_k, a_k) \\ &= \eta_k^T \begin{bmatrix} H_{\bar{x}\bar{x}} & H_{\bar{x}u} & H_{\bar{x}a} \\ H_{u\bar{x}} & H_{uu} & H_{ua} \\ H_{a\bar{x}} & H_{au} & H_{aa} \end{bmatrix} \eta_k = \eta_k^T H \eta_k \end{aligned} \quad (30)$$

where

$$\begin{aligned} H &= H^T, \eta_k = [\bar{x}_k \ u_k \ a_k]^T = [\bar{x}_k \ u_k \ a_k^1 \ \dots \ a_k^q]^T, \\ H_{\bar{x}a} &= [H_{\bar{x}a^1} \ H_{\bar{x}a^2} \ \dots \ H_{\bar{x}a^q}], H_{\bar{x}a^j} = \gamma \bar{A}^T P \bar{\Gamma}^j, \\ H_{ua} &= [H_{ua^1} \ H_{ua^2} \ \dots \ H_{ua^q}], H_{ua^j} = \gamma \bar{B}^T P \bar{\Gamma}^j, \\ H_{aa} &= \begin{bmatrix} H_{a^1 a^1} & \dots & H_{a^1 a^q} \\ \vdots & \ddots & \vdots \\ H_{a^q a^1} & \dots & H_{a^q a^q} \end{bmatrix}, \\ H_{a^i a^j} &= \begin{cases} \gamma (\bar{\Gamma}^i)^T P (\bar{\Gamma}^j), & i \neq j \\ \gamma (\bar{\Gamma}^j)^T P (\bar{\Gamma}^j) - \vartheta^j \otimes I, & i = j \end{cases} \end{aligned}$$

The optimal policies for both the defender and attackers can be obtained by solving the following equations for  $u_k$  and  $a_k^j$ .

$$\frac{\partial Q(\bar{x}_k, u_k, a_k^1, \dots, a_k^q)}{\partial u_k} = 0 \quad (31)$$

$$\frac{\partial Q(\bar{x}_k, u_k, a_k^1, \dots, a_k^q)}{\partial a_k^j} = 0, \quad j \in \mathcal{Q} \quad (32)$$

which yields

$$\begin{aligned} u_k^* &= (H_{uu} - H_{ua} H_{aa}^{-1} H_{au})^{-1} \\ &\quad \times (H_{ua} H_{aa}^{-1} H_{a\bar{x}} - H_{u\bar{x}}) \bar{x}_k \end{aligned} \quad (33)$$

and

$$\begin{aligned} a_k^* &= (H_{aa} - H_{au} H_{uu}^{-1} H_{ua})^{-1} \\ &\quad \times (H_{au} H_{uu}^{-1} H_{u\bar{x}} - H_{a\bar{x}}) \bar{x}_k \end{aligned} \quad (34)$$

where  $a_k^* = [(a_k^{1,*})^T (a_k^{2,*})^T \dots (a_k^{q,*})^T]^T$ .

Next, the aim is to express the Q-function given in (30) in terms of input-output data instead of the system state  $\bar{x}_k$ .

*Lemma 2.* Assume the system (2) and (8) is observable. Then, the system state can be re-written in terms of measured input/output sequences as

$$\begin{aligned} \bar{x}_k &= M_e \bar{e}_{k-1, k-N} + M_u \bar{u}_{k-1, k-N} \\ &\quad + \sum_{j=1}^q M_a^j \bar{a}_{k-1, k-N}^j \end{aligned} \quad (35)$$

where

$$\begin{aligned}
\bar{e}_{k-1,k-N} &= [e_{k-1}^T \ e_{k-2}^T \ \cdots \ e_{k-N}^T]^T, \\
\bar{u}_{k-1,k-N} &= [u_{k-1}^T \ u_{k-2}^T \ \cdots \ u_{k-N}^T]^T, \\
\bar{a}_{k-1,k-N}^j &= [(a_{k-1}^j)^T \ (a_{k-2}^j)^T \ \cdots \ (a_{k-N}^j)^T]^T, \\
M_e &= \bar{A}^N (V_N^T V_N)^{-1} V_N^T, \\
M_u &= U_N^1 - M_e W_N^1, \quad M_a^{j-1} = U_N^j - M_e W_N^j, \\
V_N &= [(\bar{C}\bar{A}^{N-1})^T \ (\bar{C}\bar{A}^{N-2})^T \ \cdots \ (\bar{C}\bar{A})^T \ \bar{C}^T]^T, \\
U_N^1 &= [\bar{B} \ \bar{A}\bar{B} \ \cdots \ \bar{A}^{N-1}\bar{B}], \\
U_N^j &= [\bar{\Gamma}^{j-1} \ \bar{A}\bar{\Gamma}^{j-1} \ \cdots \ \bar{A}^{N-1}\bar{\Gamma}^{j-1}], \\
W_N^1 &= \begin{bmatrix} 0 & \bar{C}\bar{B} & \bar{C}\bar{A}\bar{B} & \cdots & \bar{C}\bar{A}^{N-2}\bar{B} \\ 0 & 0 & \bar{C}\bar{B} & \cdots & \bar{C}\bar{A}^{N-3}\bar{B} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \bar{C}\bar{B} \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}, \\
W_N^j &= \begin{bmatrix} 0 & \bar{C}\bar{\Gamma}^{j-1} & \bar{C}\bar{A}\bar{\Gamma}^{j-1} & \cdots & \bar{C}\bar{A}^{N-2}\bar{\Gamma}^{j-1} \\ 0 & 0 & \bar{C}\bar{\Gamma}^{j-1} & \cdots & \bar{C}\bar{A}^{N-3}\bar{\Gamma}^{j-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \bar{C}\bar{\Gamma}^{j-1} \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}, \\
j &= 2, 3, \dots, q+1.
\end{aligned}$$

**Proof.** The proof can be easily completed and omitted here due to the space limitation.

It should be pointed out that there exists a constant  $\kappa$  such that  $\text{rank}(V_N) < n+p$  for  $N < \kappa$  and that  $\text{rank}(V_N) = n+p$  for  $N \geq \kappa$ , where  $\kappa$  is the observability index. Consequently, one can choose  $N \geq \kappa$ , and  $V_N$  has full column rank  $n+p$ . In addition,  $\bar{e}_{k-1,k-N}$ ,  $\bar{u}_{k-1,k-N}$  and  $\bar{a}_{k-1,k-N}^j$  are the tracking error output, control input and  $j^{\text{th}}$  false-data input, respectively.  $V_N$  is the observability matrix, and  $U_N^i$ ,  $i = 1, 2, \dots, q+1$  are the controllability matrices with respect to control and attacks.  $W_N^i$ ,  $i = 1, 2, \dots, q+1$  are the Toeplitz matrices. It can be seen from Lemma 1 that the augmented system state  $\bar{x}_k$  is represented as the measured input-output sequences. Next, the Q-function (30) will be described by the same input-output sequences.

Define the following vector

$$\xi_k = \begin{bmatrix} \bar{e}_{k-1,k-N} \\ \bar{u}_{k-1,k-N} \\ \bar{a}_{k-1,k-N} \\ u_k \\ a_k \end{bmatrix} \in \mathbb{R}^{[(m+1)p+m]N+(p+1)m} \quad (36)$$

Then, the Q-function can be re-written in the form of (37), where  $M_a = [M_a^1 \ M_a^2 \ \cdots \ M_a^q]$ . As a result, the optimal defender's policy  $u_k^*$  and attackers' policy  $a_k^*$  can be obtained by solving (31) and (32) simultaneously, where  $Q(\bar{x}_k, u_k, a_k) = Q(\xi_k)$ , which yields

$$\begin{aligned}
u_k^* &= (\bar{H}_{uu} - \bar{H}_{ua}(\bar{H}_{aa})^{-1}\bar{H}_{au})^{-1} \\
&\quad \times (\bar{H}_{ua}(\bar{H}_{aa})^{-1}\phi_{k-1,k-N} - \psi_{k-1,k-N}) \quad (38)
\end{aligned}$$

and

$$\begin{aligned}
a_k^* &= (\bar{H}_{aa} - \bar{H}_{au}(\bar{H}_{uu})^{-1}\bar{H}_{ua})^{-1} \\
&\quad \times (\bar{H}_{au}(\bar{H}_{uu})^{-1}\psi_{k-1,k-N} - \phi_{k-1,k-N}) \quad (39)
\end{aligned}$$

where

$$\begin{aligned}
\phi_{k-1,k-N} &= \bar{H}_{ae}\bar{e}_{k-1,k-N} + \bar{H}_{au}\bar{u}_{k-1,k-N} \\
&\quad + \bar{H}_{aa}\bar{a}_{k-1,k-N}, \\
\psi_{k-1,k-N} &= \bar{H}_{ue}\bar{e}_{k-1,k-N} + \bar{H}_{uu}\bar{u}_{k-1,k-N} \\
&\quad + \bar{H}_{ua}\bar{a}_{k-1,k-N}, \\
\bar{a}_{k-1,k-N} &= [(\bar{a}_{k-1,k-N}^1)^T \ \cdots \ (\bar{a}_{k-1,k-N}^q)^T]^T.
\end{aligned}$$

The following theorem can be obtained immediately from the above analysis, and the proof is omitted here.

*Theorem 3.* Assume the augmented systems (8)-(9) are observable. The optimal policies for both defender and attackers can be calculated according to (38) and (39), respectively, which are functions of measured input/output sequences and independent of the system state.

By combining (29) and (37), the input/output form of Bellman equation for Q-function can be re-written as

$$\begin{aligned}
\xi_k^T \bar{H} \xi_k &= \bar{x}_k^T Q_x \bar{x}_k + u_k^T R u_k - a_k^T (\vartheta \otimes I) a_k \\
&\quad + \gamma \xi_{k+1}^T \bar{H} \xi_{k+1} \quad (40)
\end{aligned}$$

where  $u_{k+1}$  and  $a_{k+1}$  can be calculated by

$$\begin{aligned}
u_{k+1} &= (\bar{H}_{uu} - \bar{H}_{ua}(\bar{H}_{aa})^{-1}\bar{H}_{au})^{-1} \\
&\quad \times (\bar{H}_{ua}(\bar{H}_{aa})^{-1}\phi_{k,k-N+1} - \psi_{k,k-N+1})
\end{aligned}$$

and

$$\begin{aligned}
a_{k+1} &= (\bar{H}_{aa} - \bar{H}_{au}(\bar{H}_{uu})^{-1}\bar{H}_{ua})^{-1} \\
&\quad \times (\bar{H}_{au}(\bar{H}_{uu})^{-1}\psi_{k,k-N+1} - \phi_{k,k-N+1})
\end{aligned}$$

Now, we linearly parameterize the Q-function as follows

$$\begin{aligned}
\xi_k^T \bar{H} \xi_k &= \bar{H}_{11}(\xi_k^1)^2 + 2\bar{H}_{12}\xi_k^1 \xi_k^2 + \cdots + 2\bar{H}_{1\ell}\xi_k^1 \xi_k^\ell \\
&\quad + \bar{H}_{22}(\xi_k^2)^2 + 2\bar{H}_{23}\xi_k^2 \xi_k^3 + \cdots + 2\bar{H}_{2\ell}\xi_k^2 \xi_k^\ell \\
&\quad + \cdots + \bar{H}_{\ell\ell}(\xi_k^\ell)^2 \\
&= h^T \bar{\xi}_k \quad (41)
\end{aligned}$$

where

$$\begin{aligned}
\bar{H} &= \bar{H}^T, \quad \xi_k = [\xi_k^1 \ \xi_k^2 \ \cdots \ \xi_k^\ell]^T, \\
\ell &= [(m+1)p+m]N + (p+1)m, \\
h &= [\bar{H}_{11} \ 2\bar{H}_{12} \ \cdots \ 2\bar{H}_{1\ell} \ \bar{H}_{22} \ \cdots \ 2\bar{H}_{2\ell} \ \cdots \ \bar{H}_{\ell\ell}]^T, \\
\bar{\xi}_k &= [(\xi_k^1)^2 \ \xi_k^1 \xi_k^2 \ \cdots \ \xi_k^1 \xi_k^\ell \ (\xi_k^2)^2 \ \cdots \ \xi_k^2 \xi_k^\ell \ \cdots \ (\xi_k^\ell)^2]^T.
\end{aligned}$$

It should be emphasized that the unknown matrix  $\bar{H} \in \mathbb{R}^{\ell \times \ell}$  has  $\frac{1}{2}\ell(\ell+1)$  unknown elements due to  $\bar{H}_{ij} = \bar{H}_{ji}$ . From (40) and (41), one can obtain that

$$h^T \bar{\xi}_k = \bar{x}_k^T Q_x \bar{x}_k + u_k^T R u_k - a_k^T (\vartheta \otimes I) a_k + h^T \bar{\xi}_{k+1} \quad (42)$$

which is a key equation in the following Q-learning algorithms. Now, we are ready to use the Q-learning approach to learn the corresponding Q-function matrix  $\bar{H}$ . Policy iteration and value iteration algorithms using Q-learning technique are given in Algorithms 1 and 2, respectively. In Algorithms 1 and 2,  $a_k^j$ ,  $\phi_{k-1,k-N}^{j+1}$  and  $\psi_{k-1,k-N}^{j+1}$  are defined as follows

$$\begin{aligned}
a_k^j &= [(a_k^{1,j})^T, (a_k^{2,j})^T, \dots, (a_k^{q,j})^T]^T \\
\phi_{k-1,k-N}^{j+1} &= \bar{H}_{ae}^{j+1} \bar{e}_{k-1,k-N} + \bar{H}_{au}^{j+1} \bar{u}_{k-1,k-N} \\
&\quad + \bar{H}_{aa}^{j+1} \bar{a}_{k-1,k-N} \\
\psi_{k-1,k-N}^{j+1} &= \bar{H}_{ue}^{j+1} \bar{e}_{k-1,k-N} + \bar{H}_{uu}^{j+1} \bar{u}_{k-1,k-N} \\
&\quad + \bar{H}_{ua}^{j+1} \bar{a}_{k-1,k-N}
\end{aligned}$$

$$Q(\xi_k) = \xi_k^T \begin{bmatrix} M_e^T H_{\bar{x}\bar{x}} M_e & M_e^T H_{\bar{x}\bar{x}} M_u & M_e^T H_{\bar{x}\bar{x}} M_a & M_e^T H_{\bar{x}u} & M_e^T H_{\bar{x}a} \\ M_u^T H_{\bar{x}\bar{x}} M_e & M_u^T H_{\bar{x}\bar{x}} M_u & M_u^T H_{\bar{x}\bar{x}} M_a & M_u^T H_{\bar{x}u} & M_u^T H_{\bar{x}a} \\ M_a^T H_{\bar{x}\bar{x}} M_e & M_a^T H_{\bar{x}\bar{x}} M_u & M_a^T H_{\bar{x}\bar{x}} M_a & M_a^T H_{\bar{x}u} & M_a^T H_{\bar{x}a} \\ H_{u\bar{x}} M_e & H_{u\bar{x}} M_u & H_{u\bar{x}} M_a & H_{uu} & H_{ua} \\ H_{a\bar{x}} M_e & H_{a\bar{x}} M_u & H_{a\bar{x}} M_a & H_{au} & H_{aa} \end{bmatrix} \xi_k = \xi_k^T \begin{bmatrix} \bar{H}_{\bar{e}\bar{e}} & \bar{H}_{\bar{e}\bar{u}} & \bar{H}_{\bar{e}\bar{a}} & \bar{H}_{\bar{e}u} & \bar{H}_{\bar{e}a} \\ \bar{H}_{\bar{u}\bar{e}} & \bar{H}_{\bar{u}\bar{u}} & \bar{H}_{\bar{u}\bar{a}} & \bar{H}_{\bar{u}u} & \bar{H}_{\bar{u}a} \\ \bar{H}_{\bar{a}\bar{e}} & \bar{H}_{\bar{a}\bar{u}} & \bar{H}_{\bar{a}\bar{a}} & \bar{H}_{\bar{a}u} & \bar{H}_{\bar{a}a} \\ \bar{H}_{u\bar{e}} & \bar{H}_{u\bar{u}} & \bar{H}_{u\bar{a}} & \bar{H}_{uu} & \bar{H}_{ua} \\ \bar{H}_{a\bar{e}} & \bar{H}_{a\bar{u}} & \bar{H}_{a\bar{a}} & \bar{H}_{au} & \bar{H}_{aa} \end{bmatrix} \xi_k \quad (37)$$

---

**Algorithm 1 : Policy Iteration Using Q-Learning**

---

1. **Initialization** : Set  $j = 0$ , select stabilizing defender's policy  $u_k^0$  and attackers' policy  $a_k^0$ , and choose  $H^0 = (H^0)^T$

2. **Policy Evaluation** : Solve for  $h^{j+1}$   
 $-(h^{j+1})^T \bar{\xi}_k + \bar{x}_k^T Q_x \bar{x}_k + (u_k^j)^T R(u_k^j)$   
 $-(a_k^j)^T (\vartheta \otimes I)(a_k^j) + \gamma(h^{j+1})^T \bar{\xi}_{k+1} = 0$

3. **Policy Improvement** :

$$u_k^{j+1} = (\bar{H}_{uu}^{j+1} - \bar{H}_{ua}^{j+1}(\bar{H}_{aa}^{j+1})^{-1}\bar{H}_{au}^{j+1})^{-1} \\ \times (\bar{H}_{ua}^{j+1}(\bar{H}_{aa}^{j+1})^{-1}\phi_{k-1,k-N}^{j+1} - \psi_{k-1,k-N}^{j+1}) \\ a_k^{j+1} = (\bar{H}_{aa}^{j+1} - \bar{H}_{au}^{j+1}(\bar{H}_{uu}^{j+1})^{-1}\bar{H}_{ua}^{j+1})^{-1} \\ \times (\bar{H}_{au}^{j+1}(\bar{H}_{uu}^{j+1})^{-1}\psi_{k-1,k-N}^{j+1} - \phi_{k-1,k-N}^{j+1})$$

4. **Stop if**

$$\|\bar{H}^{j+1} - \bar{H}^j\| < \epsilon$$


---

**Algorithm 2 : Value Iteration Using Q-Learning**

---

1. **Initialization** : Set  $j = 0$ , select any defender's policy  $u_k^0$  and attackers' policy  $a_k^0$ , and choose  $H^0 = (H^0)^T$

2. **Policy Evaluation** : Solve for  $h^{j+1}$   
 $(h^{j+1})^T \bar{\xi}_k = \bar{x}_k^T Q_x \bar{x}_k + (u_k^j)^T R(u_k^j)$   
 $-(a_k^j)^T (\vartheta \otimes I)(a_k^j) + \gamma(h^j)^T \bar{\xi}_{k+1}$

3. **Policy Improvement** :

$$u_k^{j+1} = (\bar{H}_{uu}^{j+1} - \bar{H}_{ua}^{j+1}(\bar{H}_{aa}^{j+1})^{-1}\bar{H}_{au}^{j+1})^{-1} \\ \times (\bar{H}_{ua}^{j+1}(\bar{H}_{aa}^{j+1})^{-1}\phi_{k-1,k-N}^{j+1} - \psi_{k-1,k-N}^{j+1}) \\ a_k^{j+1} = (\bar{H}_{aa}^{j+1} - \bar{H}_{au}^{j+1}(\bar{H}_{uu}^{j+1})^{-1}\bar{H}_{ua}^{j+1})^{-1} \\ \times (\bar{H}_{au}^{j+1}(\bar{H}_{uu}^{j+1})^{-1}\psi_{k-1,k-N}^{j+1} - \phi_{k-1,k-N}^{j+1})$$

4. **Stop if**

$$\|\bar{H}^{j+1} - \bar{H}^j\| < \epsilon$$


---

In order to solve  $\bar{H}^{j+1}$  by using (40) recursively, the number of samples  $\bar{\xi}_k$  should satisfy  $v \geq \frac{1}{2}\ell(\ell + 1)$ . Let  $\Xi = [\bar{\xi}_k^1, \bar{\xi}_k^2, \dots, \bar{\xi}_k^v]$  and  $\Lambda = [\rho_k^1, \rho_k^2, \dots, \rho_k^v]$ , where  $\rho_k^i = [\bar{x}_k^T Q_x \bar{x}_k + (u_k^j)^T R(u_k^j) - \vartheta(a_k^j)^T (a_k^j)]^i + \gamma(h^j)^T \bar{\xi}_{k+1}^i$ ,  $i = 1, 2, \dots, v$ . Then, one can obtain that  $(h^{j+1})^T \Xi = \Lambda$ , which results in

$$h^{j+1} = (\Xi \Xi^T)^{-1} \Xi \Lambda^T \quad (43)$$

It is noted that, however, the defender's policy  $u_k$  and the attacker's policy  $a_k$  are dependent on the measured input-output  $e_{k-1,k-N}$ ,  $u_{k-1,k-N}$  and  $a_{k-1,k-N}$ . Consequently, the matrix  $\Xi \Xi^T$  is not invertible. To address this issue, probing noise should be added into the system dynamics, and the following persistence of excitation (PE) condition must be satisfied.

*Definition 4.* A  $q$ -vector sequence  $h = [h_1, h_2, \dots, h_q]^T$  is said to be persistently exciting over an interval  $[k+1, k+l]$  if for some constant  $\delta > 0$

$$\sum_{i=k+1}^{k+l} h_i (h_i)^T \geq \delta I \quad (44)$$

In order to satisfy the persistence of excitation condition, the actual defender's policy and attacker's policy are generated by

$$\hat{u}_k = u_k + n_k^1, \quad \hat{a}_k^j = a_k^j + n_k^{2,j} \quad (45)$$

where  $n_k^1$  and  $n_k^{2,j}$ ,  $j = 1, 2, \dots, q$ , are probing noise signals.

*Remark 5.* If the state almost converges to the desired position and becomes stationary, then the persistence of excitation condition is no longer satisfied. An exploratory signal consisting of sinusoids of varying frequencies can be added to the policies of the defender and attackers to ensure PE qualitatively. Consequently,  $n_k^1$  and  $n_k^{2,j}$  in (45) can be chosen as sinusoids of varying frequencies, exponentially decaying noise, or Gaussian noise.

According to the method given in Rizvi et al. (2018), the following theorem can be easily obtained, and the corresponding proof is omitted here.

*Theorem 6.* Assume the system is controllable and observable. The sequences  $u_k^j$  and  $a_k^{i,j}$ ,  $i \in \mathcal{Q}$ ,  $j = 1, 2, \dots, \infty$ , generated by Q-learning Algorithm 7 or Algorithm 8, can converge to the optimal strategies (38) and (39) for the defender and the attackers when the system is sufficiently excited.

#### 4. AN ILLUSTRATIVE EXAMPLE

In this section, an example is given to demonstrate the main results proposed in this paper.

A linear discrete-time system is given as

$$x_{k+1} = Ax_k + Bu_k + B(\Gamma^1 a_k^1 + \Gamma^2 a_k^2) \\ y_k = Cx_k$$

with

$$A = \begin{bmatrix} -1.5 & 0.3 & 1 \\ 1.1 & 0.7 & -0.5 \\ 0.5 & -0.2 & 1.9 \end{bmatrix}, \quad B = \begin{bmatrix} 0.2 & -0.3 \\ 0.6 & -1 \\ -1.1 & 0.8 \end{bmatrix}, \\ \Gamma^1 = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad \Gamma^2 = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}, \quad C = [1 \quad 1 \quad 1].$$

Note that the eigenvalues of system matrix  $A$  are  $\lambda_1 = -1.8021$ ,  $\lambda_2 = 0.8380$  and  $\lambda_3 = 2.0641$ . Therefore, the open-loop system is unstable. The reference trajectory is generated by

$$y_{k+1}^r = -y_k^r$$

Other parameters are given as follow.

$$\gamma = 0.88, \quad \theta_1 = 0.83, \quad \theta_2 = 0.81, \quad \epsilon = 10^{-6}, \\ Q_e = 0.9, \quad R = \begin{bmatrix} 0.25 & 0 \\ 0 & 0.001 \end{bmatrix}.$$

The simulation results can be obtained according to Algorithm 2 and are depicted in Fig.1 - Fig.4, which show that the tracking error can converge to zero, i.e., the system output  $y_k$  can track the reference input  $y_k^r$ .

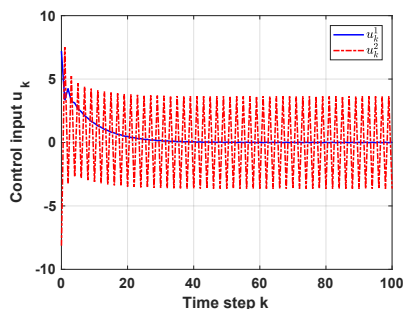


Fig. 1. Control input  $u_k$

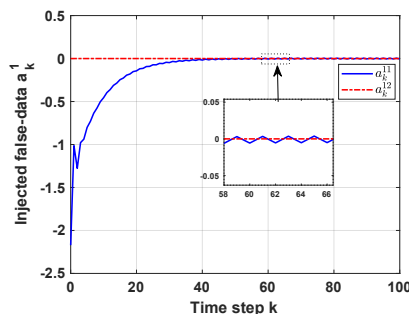


Fig. 2. False-data injected by attacker 1

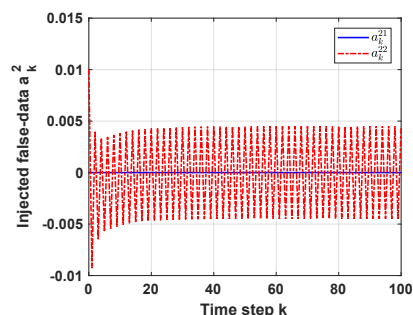


Fig. 3. False-data injected by attacker 2

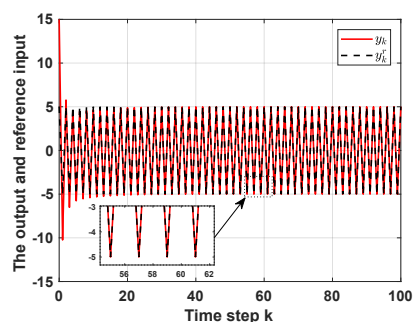


Fig. 4. System output and the reference trajectory

## 5. CONCLUSION

In this paper, game theory was used to investigate the optimal tracking control problem in the presence of false-data-injection attacks. Then, the Q-learning method was developed to solve the GARE online without requiring the knowledge of augmented system dynamics. Moreover, the Q-function was expressed in terms of only measured input-output data, and the policies for both sides were

generated by Q-learning algorithm, where the system was assumed to be sufficiently excited. The simulation results have shown that the system output can track the given reference trajectory under FDI attacks.

## REFERENCES

- A. Teixeira, I. Shames, H. Sandberg, and K. H. Johansson. A secure control framework for resource-limited adversaries, *Automatica*, 51, 135–148, 2015.
- D. Persis, and P. Tesi. Networked control of nonlinear systems under denial-of-service, *Syst. Contr. Lett.*, 96, 124–131, 2016.
- J. H. Qin, M. L. Li, L. Shi, and X. H. Yu. Optimal denial-of-service attack scheduling with energy constraint over packet-dropping networks, *IEEE Trans. Autom. Control*, 63(6), 1648–1663, 2018.
- K. M. Ding, Y. Z. Li, D. E. Quevedo, S. Dey, and L. Shi. A multi-channel transmission schedule for remote state estimation under DoS attacks, *Automatica*, 78, 194–201, 2017.
- E. Kung, S. Dey, and L. Shi. The performance and limitations of  $\epsilon$ -stealthy attacks on higher order systems, *IEEE Trans. Autom. Control*, 62(2), 941–947, 2017.
- C. Z. Bai, V. Gupta, and F. Pasqualetti. On Kalman filtering with compromised sensors: Attacks stealthiness and performance bounds, *IEEE Trans. Autom. Control*, 62(12), 6641–6648, 2017.
- L. Hu, Z. D. Wang, Q. L. Han, and X. H. Liu. State estimation under false data injection attacks: Security analysis and system protection, *Automatica*, 87, 176–183, 2018.
- Y. L. Mo, and B. Sinopoli. Secure control against replay attacks, *In Proc. 47 Annual Allerton Conference*, pp. 911–918, 2009.
- M. H. Zhu, and S. Martinez. On the performance analysis of resilient networked control systems under replay attacks, *IEEE Trans. Autom. Control*, 59(3), 804–808, 2014.
- B. Chen, D. W. C. Ho, G. Hu, and L. Yu. Secure fusion estimation for bandwidth constrained cyber-physical systems under replay attacks, *IEEE Trans. Cybern.*, 48(6), 1862–1876, 2018.
- C. Z. Bai, F. Pasqualetti, and V. Gupta. Data-injection attacks in stochastic control systems: Detectability and performance tradeoffs, *Automatica*, 82, 251–260, 2017.
- Y. Z. Li, D. W. Shi, and T. W. Chen. False data injection attacks on networked control systems: A Stackelberg-game analysis, *IEEE Trans. Autom. Control*, 63(10), 3503–3509, 2018.
- S. A. A. Rizvi, and Z. L. Lin. Output feedback Q-learning for discrete-time linear zero-sum games with application to the H-infinity control, *Automatica*, 95, 213–221, 2018.
- D. Ye, S. P. Luo. A co-design methodology for cyber-physical systems under actuator fault and cyber attack, *J. Franklin Inst.*, 356, 1856–1879, 2019.
- D. R. Ding, Z. D. Wang, D. W. C. Ho, and G. L. Wei. Observer-based event triggering consensus control for multi-agent systems with lossy sensors and cyber attacks, *IEEE Trans. Cybern.*, 47(8), 1936–1947, 2017.