

State estimation under false data injection attacks: Security analysis and system protection

Hu, L., Wang, Z., Han, Q., & Liu, X. (2018). State estimation under false data injection attacks: Security analysis and system protection. *Automatica*, *87*, 176-183. https://doi.org/10.1016/j.automatica.2017.09.028

Published in: Automatica

Document Version: Peer reviewed version

Queen's University Belfast - Research Portal: Link to publication record in Queen's University Belfast Research Portal

Publisher rights

© 2017 Elsevier Ltd. This manuscript version is made available under the CC-BY-NC-ND 4.0 license http://creativecommons.org/licenses/bync-nd/4.0/,which permits distribution and reproduction for noncommercial purposes, provided the author and source are cited.

General rights

Copyright for the publications made accessible via the Queen's University Belfast Research Portal is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The Research Portal is Queen's institutional repository that provides access to Queen's research output. Every effort has been made to ensure that content in the Research Portal does not infringe any person's rights, or applicable UK laws. If you discover content in the Research Portal that you believe breaches copyright or violates any law, please contact openaccess@qub.ac.uk.

Open Access

This research has been made openly available by Queen's academics and its Open Research team. We would love to hear how access to this research benefits you. – Share your feedback with us: http://go.qub.ac.uk/oa-feedback

State Estimation under False Data Injection Attacks: Security Analysis and System Protection

Liang Hu, Zidong Wang and Xiaohui Liu

Abstract

In this paper, the security issue in the state estimation problem is investigated for networked control systems. The communication channels between the sensors and the remote estimator are vulnerable to attacks from malicious adversaries. The false data injection attacks (FDIAs) are considered. We aim to find the so-called *insecurity* conditions under which the estimation system is insecure in the sense that there exist FDIAs that can bypass the anomaly detector but still lead to unbounded estimation errors. In particular, a *new* necessary and sufficient condition for the insecurity is derived in the case that all communication channels are compromised by the adversary. Furthermore, a specific algorithm is proposed for generating attacks with which the estimation system is insecure. Moreover, for the insecure system, we propose a system protection scheme through which only a few (rather than all) communication channels require protection against FDIAs. A simulation example is utilized to demonstrate the usefulness of the proposed conditions/algorithms in the secure estimation problem for a flight vehicle.

Index Terms

False data injection attacks; State estimation, Security analysis, Networked control systems

I. INTRODUCTION

The first-ever cyber-attack in real-world control systems was reported in 2010 [4]. Specifically, the Stuxnet worms were injected in the programming logic controllers leading to subsequent modifications of the control systems and eventual damages of nearly 1000 centrifuges. Since then, the cybersecurity of NCSs has been a hot topic of research that stirs considerable interest. In general, two kinds of attacks have been studied in NCSs [23], one is the denial-of-service (DoS) attack that violates data *availability* through blocking information flows between different components of NCSs, and the other is the deception attack that violates data *integrity* through modifying the data packets. Compared with DoS attacks, deception attacks are more difficult to detect because the adversary could keep the deception attacks stealthy to the anomaly detector in NCSs through deliberately designing the attack sequences.

The deception attacks have been first considered in [12] for the state estimation problems of power systems modelled by *static* system models. Some attack detection methods have been proposed in [7], [10] from the defenders' perspective, and the minimum number of comprised sensors that needed to launch deception attacks has been investigated in [22] from attackers' points of view. As for *dynamic* systems, when the system model is unknown to the adversary, a specific type of deception attack called replay attack has been proposed and analysed in [15], [18], [25]. In the case that the dynamic model is known to the adversary, another type of deception attack, namely, false data injection attack (FDIA), has recently been put forward. For deterministic systems without stochastic noises, fundamental issues such as detectability and identifiability for FDIAs have been analysed in [5], [6], [20]

This work was supported in part by the Engineering and Physical Sciences Research Council (EPSRC) of the U.K., the Royal Society of the U.K., and the Alexander von Humboldt Foundation of Germany.

L. Hu, Z. Wang and X. Liu are with the Department of Computer Science, Brunel University London, Uxbridge, Middlesex, UB8 3PH, United Kingdom. (Email: Zidong.Wang@brunel.ac.uk)

and efficient control/estimation algorithms have been developed against FDIAs [21]. In [19], the data encryption scheme (together with time-stamp techniques) has been adopted to detect the deception attacks and compensate the side-effects.

As is well known, stochastic models have come to play a more and more important role in characterizing noisy phenomena from real-world systems. Accordingly, it is of practical significance to investigate the cybersecurity of stochastic dynamic systems. As pointed in [9], the detection task of malicious behaviours for stochastic systems (with external noises) is more difficult than that for deterministic (without stochastic noises) because of the fact that the injected attack by the adversary could be mistaken as a type of noises by the protection devices. A secure state estimation algorithm has been proposed in [13] for stochastic dynamic systems where a key assumption of sparse observability has been made. This assumption implies that only a part of (rather than all) sensors are attacked and the system is still observable using the set of unattacked sensors. While the results reported in [13] are indeed interesting, it is quite possible that the adversary attacks at a large number of (or even all) sensors, in which case the system cannot be guaranteed to be "sparsely observable". Motivated by the above observation, we aim to examine the system vulnerability under cyber-attacks without the assumption of sparse observability. More specifically, we investigate the case where the attacker could inject false data into measurements from any sensor and, accordingly, the main results obtained would constitute one of the main contributions of our paper.

In this paper, we focus on the remote state estimation problem for a class of *stochastic* systems under possible FDIA attacks where a χ^2 detector is employed to monitor the state estimates. Note that FDIAs have been considered in [11], [14], [16], [17] for state estimation problems, and in [24] for distributed control problems of stochastic systems equipped with χ^2 detectors. In particular, an approximation method has been proposed in [16], [17] to analyse the cybersecurity of the system by calculating the estimation error bound caused by the FDIA attacks, and some *insecurity* conditions have been derived in [11], [14] to determine whether or not there exists FDIA attacks which can cause unbounded estimation error for the state estimation system. Nevertheless, a thorough investigation reveals that 1) there is still room to improve the existing insecurity conditions; and 2) there is also an engineering need to develop system protection scheme by using only necessary number of communication channels requiring protection against FDIAs.

In this paper, we aim to propose new insecurity conditions for state estimation problems under FDIAs. Specifically, for the case when all communication channels are compromised by the adversary, we propose a *new* necessary and sufficient condition under which the system is insecure in the sense that the estimation error caused by FDIAs is unbounded. Such a new condition is shown to be concise that simplifies the existing results. For the case when only parts of the communication channels are compromised by the adversary, a sufficient condition is proposed as well. Furthermore, to protect the overall NCSs from FDIAs, we propose a criterion which determines a sufficient number of communication channels that require protection. According to the criterion, only necessary number of (rather than all) communication channels need to be protected in order to make the overall system secure against the FDIAs. The contributions of the paper are summarized as follows: *1) new security criteria are proposed for state estimation systems under FDIAs and, in the case that all communication channels are compromised by the adversary, our criteria are shown to be necessary and sufficient that simplify the existing ones; <i>2) an effective protection scheme is proposed for the system which is insecure under FDIAs; and 3) the developed criteria are applied to security analysis and system protection in the state estimation system of a flight vehicle.*

The remainder of this paper is organized as follows. The security problem of state estimation system under cyber-attacks are formulated in Section II. In Section III, we analyse the system security under FDIAs for two cases and further propose the system protection scheme. Examples for illustration are given in Section IV and we conclude the paper in Section V.



Fig. 1. Diagram of state estimation problem under cyber-attacks

Notation: \mathbb{N} , \mathbb{R} and \mathbb{C} denote, respectively, the set of non-negative integers, the set of all real numbers, and the set of all complex numbers. $\{x(k)\}$ denotes an infinite sequence $x(1), x(2), \dots, x(k), \dots$. $\mathbb{R}^{n \times m}$ ($\mathbb{C}^{n \times m}$) denotes the set of all $n \times m$ real (complex) matrices, and \mathbb{R}^n denotes the n dimensional Euclidean space. For $\alpha \in \mathbb{C}$, $\operatorname{Re}(\alpha)$ and $|\alpha|$ denote its real part and its modulus, respectively. For $a \in \mathbb{R}^n$, ||a|| denotes its l_2 norm. For a matrix $P \in \mathbb{R}^{n \times m}$, P^T , P^{-1} , $\operatorname{Tr}\{P\}$ and $\operatorname{Rk}\{P\}$ represent its transpose, inverse, trace, and rank, respectively. For square matrix A, $\det(A)$ stands for the determinant of A, and $\rho(A)$ stands for the spectral radius of A. $\operatorname{diag}\{\dots\}$ and I denotes the $m \times m$ -dimensional identity (zeros) matrix. I_m^s denotes the s-th column of $m \times m$ -dimensional identity matrix I_m , e.g., $I_m^s = [\overbrace{0,\dots,0,1,0,\dots,0}^{s-1}]^T$.

II. PROBLEM FORMULATION

In this section, we describe the model of false data injection attack (FDIA) and analyse how the injected attacks affect the estimation system. The structure of the state estimation system under cyber-attacks is shown in Fig. 1. For presentation convenience, we first introduce the estimation system without cyber-attacks (i.e., $y^a(k) = y(k)$ in Fig. 1).

A. State estimation without cyber-attacks

Let the physical plant be given by:

$$\mathcal{P}: \begin{cases} x(k+1) = Ax(k) + \omega(k) \\ y(k) = Cx(k) + \nu(k) \end{cases}$$
(1)

where $x(k) \in \mathbb{R}^n$ is the system state, $y(k) = [y_1(k), \dots, y_m(k)]^T \in \mathbb{R}^m$ is the measurement output, and $y_i(k)$ is the output of the *i*th sensor (labelled as S_i in Fig. 1) at time instant k. The initial state x(0) has mean $\bar{x}(0)$ and covariance $\Sigma(0)$, the process noise $\omega(k) \in \mathbb{R}^n$ and the measurement noise $\nu(k) \in \mathbb{R}^m$ are assumed to be mutually uncorrelated zero-mean random signals with known covariance matrices W and R, respectively. It is assumed that (A, C) is observable. The following time-invariant state estimator is proposed:

$$\mathcal{E} : \begin{cases} \hat{x}(k+1) = A\hat{x}(k) + Kz(k+1) \\ z(k+1) = y(k+1) - CA\hat{x}(k) \end{cases}$$
(2)

where $\hat{x}(k+1)$ and z(k+1) are the state estimate and the estimation residual at time instant k+1, respectively. Throughout this paper, we assume that the estimator converges to its steady state.

Defining the estimation error $\tilde{x}(k+1) \triangleq x(k+1) - \hat{x}(k+1)$, the dynamics of the estimation error follows from (1) and (2) as follows:

$$\tilde{x}(k+1) = (I - KC)(A\tilde{x}(k) + \omega(k)) - K\nu(k+1).$$
(3)

It is well known that the estimator is stable if and only if the matrix (I - KC)A is stable [8]. In this paper, it is assumed that the estimator is stable by choosing appropriate estimator gain K.

Failure detectors are often used to detect abnormal operations. In this paper, we assume that a χ^2 failure detector is deployed. At each time instant k, the χ^2 failure detector first computes the value $g(k) = z^T(k)(C\Sigma C^T + R)^{-1}z(k)$ where Σ is the steady estimation error covariance, and then compares g(k) with a prescribed threshold α . If $g(k) > \alpha$, then an alarm will be triggered. When the system operates normally (i.e. without attacks), g(k) satisfies a χ^2 distribution implying low probability of a large g(k) [1].

B. False data injection attack

In this subsection, we introduce the model of false data injection attack (FDIA) and then investigate how it affects the estimation dynamics. Assume that the adversary has perfect knowledge about the system model, that is, the values of all the matrices A, C, K, W and R described in Subsection II-A are known by the attacker. We also assume that the attacker has the ability to inject false data over the communication channels between the sensors and the estimator. Under FDIAs, the measurement output received by the estimator is given as follows:

$$y^{a}(k) = Cx(k) + a(k) + \nu(k) = Cx(k) + B_{a}a^{0}(k) + \nu(k)$$
(4)

where $a(k) \in \mathbb{R}^m$ represents the false data injected by the attacker at time instant k. The attack vector is described by $a(k) = B_a a^0(k)$ where the injection matrix $B_a = \text{diag}\{\gamma_1, \ldots, \gamma_m\}$. Here, $\gamma_i = 1$ if the attacker is able to inject false data into the *i*th communication channel, otherwise $\gamma_i = 0$. The matrix B_a reflects which communication channels can be compromised by the attacker. Specifically, $B_a = 0$ means that no FDIAs can be injected into any communication channel, and $B_a = I_m$ implies that the attacker has the ability to inject FDIA into all communication channels.

With the compromised measurement $y^a(k)$, based on the estimator \mathcal{E} in (2), the dynamics of state estimation can be derived as follows:

$$\hat{x}^{a}(k+1) = A\hat{x}^{a}(k) + Kz^{a}(k+1)$$

$$z^{a}(k+1) = y^{a}(k+1) - CA\hat{x}^{a}(k)$$
(5)

where $\hat{x}^a(k+1)$ and $z^a(k+1)$ are the state estimation and the estimation residual of system (1) at time k+1 using the compromised measurement (4), respectively. Without loss of generality, we assume that the attack begins at time instant 1 and $\hat{x}^a(0) = \hat{x}(0)$.

To take into account the effect of FDIAs on the state estimation of system (1), we define the difference between the state estimates and estimation residual of system (1) (without FDIAs) and system (4) (with FDIAs) as

$$\Delta \hat{x}(k+1) \triangleq \hat{x}^a(k+1) - \hat{x}(k+1), \Delta z(k+1) \triangleq z^a(k+1) - z(k+1).$$

For convenience, we call $\Delta \hat{x}(k+1)$ and $\Delta z(k+1)$ as the state estimation difference and the estimation residual difference, respectively. The dynamics of $\Delta z(k+1)$ and $\Delta \hat{x}(k+1)$ can be derived from (2) and (5) as follows:

$$\Delta z(k+1) = -CA\Delta \hat{x}(k) + a(k+1),$$

$$\Delta \hat{x}(k+1) = A\Delta \hat{x}(k) + K\Delta z(k+1)$$
(6)

$$= (I - KC)A\Delta\hat{x}(k) + Ka(k+1)$$
(7)

where $\Delta \hat{x}(0) = \hat{x}^{a}(0) - \hat{x}(0) = 0.$

In the considered FDIA model, the purpose of the attacker is to launch a "special" FDIA sequence under which the state estimation difference $\Delta \hat{x}(k)$ will diverge to ∞ without any alarm triggered by the χ^2 detector. In other words, the attacker aims to inject false data which would largely degrade the estimation performance without being detected by the detector.

It is known from the triangular inequality $||z^a(k)|| \le ||z(k)|| + ||\Delta z(k)||$ that, if $||\Delta z(k)||$ is small, then the χ^2 detector cannot distinguish between $z^a(k)$ and z(k) with high probability. As such, to make the attack sequence stealthy, the attacker launching the FDIA should avoid causing a large change in estimation residual difference $\Delta z(k)$ [14], which means that the inequality $||\Delta z(k)|| \le M$ should hold all the time, where M represents the tolerant level of the χ^2 detector. Obviously, a smaller value of M would result in a higher probability for the corresponding attack to be undetected. We assume that M is predetermined by the attacker. On the other hand, the attacker should design the attack sequence deliberately such that the sequence $\{\Delta \hat{x}(k)\}$ becomes unbounded, i.e, $\lim_{k\to\infty} \Delta \hat{x}(k) = \infty$.

Throughout the paper, the definition on system security is given as follows.

Definition 1: The system \mathcal{P} in (1) with estimator \mathcal{E} in (2) is called *insecure* if there exists at least one FDIA sequence $\{a(k)\}$ such that the following two conditions are satisfied simultaneously:

1) for the state estimation difference $\Delta \hat{x}(k)$,

$$\lim_{k \to \infty} \|\Delta \hat{x}(k)\| \to \infty; \tag{8}$$

2) for the estimation residual difference $\Delta z(k)$,

$$\|\Delta z(k)\| \le M,\tag{9}$$

where M is a prescribed small positive constant scalar.

In case that (8)-(9) do not hold simultaneously under FDIAs (4), the system \mathcal{P} in (1) with estimator \mathcal{E} in (2) is called *secure* under FDIAs (4).

The aim of the addressed system security problem is to analyse under what conditions there exists an FDIA that is undetectable by the fault detector but drives the bias in state estimation to infinity.

III. SECURITY ANALYSIS

In this section, we investigate the security of system \mathcal{P} in (1) with estimator \mathcal{E} in (2) for the following two cases: 1) the attacker is able to inject FDIAs into all communication channels, *i.e.*, $B_a = I_m$; and 2) the attacker can inject FDIAs into only part of the communication channels, *i.e.*, $B_a \neq I_m$.

Assume that the system matrix A in (1) has p independent eigenvectors and its Jordan form J is given by

$$J = P^{-1}AP \tag{10}$$

where

$$J = \begin{bmatrix} J_1 & 0 & 0 & \dots & 0 \\ 0 & J_2 & 0 & \dots & 0 \\ 0 & 0 & J_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & J_p \end{bmatrix}, \quad J_i = \begin{bmatrix} \lambda_i & 1 & & & \\ & \lambda_i & 1 & & \\ & & \ddots & \ddots & \\ & & & \ddots & 1 \\ & & & & \lambda_i \end{bmatrix},$$

the Jordan block $J_i \in \mathbb{C}^{n_i \times n_i}$ (i = 1, ..., p) with $|\lambda_1| \ge |\lambda_2| \ge \cdots \ge |\lambda_p|$ and $\sum_{i=1}^{i=p} n_i = n$. Denote $P = \begin{bmatrix} P_1, \ldots, P_p \end{bmatrix}$ and $Q = P^{-1} = \begin{bmatrix} Q_1^T, \ldots, Q_p^T \end{bmatrix}^T$, where $P_i \in \mathbb{C}^{n \times n_i}$ and $Q_i \in \mathbb{C}^{n_i \times n}$. If $\rho(A) \ge 1$, there exists a positive integer l satisfying $1 \le l \le p$ such that the inequality $|\lambda_1| \ge \cdots \ge |\lambda_l| \ge 1$.

 $1 > |\lambda_{l+1}| \ge \cdots \ge |\lambda_p|$ is true. Furthermore, defining $\bar{l} = \sum_{i=1}^l n_i$, we have

$$A = PJQ = \begin{bmatrix} P_o & P_c \end{bmatrix} \begin{bmatrix} \Lambda_1 & 0 \\ 0 & \Lambda_2 \end{bmatrix} \begin{bmatrix} Q_o \\ Q_c \end{bmatrix},$$
(11)

where block matrices $\Lambda_1 = \operatorname{diag}\{J_1, \dots, J_l\} \in \mathbb{C}^{\bar{l} \times \bar{l}}, \Lambda_2 = \operatorname{diag}\{J_{l+1}, \dots, J_p\} \in \mathbb{C}^{(n-\bar{l}) \times (n-\bar{l})}, P_o = [P_1, \dots, P_l],$ $P_c = [P_{l+1}, \dots, P_p], Q_o = [Q_1^T, \dots, Q_l^T]^T \text{ and } Q_c = [Q_{l+1}^T, \dots, Q_p^T]^T \text{ are of appropriate dimensions.}$

A. Case 1: $B_a = I_m$

To introduce our main results, we need the following lemmas.

Lemma 1: [3] For two matrices $M, N \in \mathbb{C}^{n \times n}$, det(MN) = det(M)det(N). Moreover, matrices MN and NMhave the same non-zero eigenvalues.

Lemma 2: For the system (1) with estimator (2), if $\rho(A) \ge 1$, the following matrix equation

$$P_c X = K \tag{12}$$

has no solution, where matrix K is the estimator gain of state estimator (2) and matrix P_c is given in (11).

Proof: It is known from Lemma 1 that the matrices (I - KC)A and A(I - KC) have the same eigenvalues. Then, it follows from $\rho((I - KC)A) < 1$ that the inequality $\rho(A(I - KC)) < 1$ holds.

Let us prove the lemma by contradiction. Assume that there exists a matrix solution \hat{X} to equation (12), then we have

$$A(I - KC) = \begin{bmatrix} P_o & P_c \end{bmatrix} \begin{bmatrix} \Lambda_1 & 0 \\ 0 & \Lambda_2 \end{bmatrix} \begin{bmatrix} Q_o \\ Q_c \end{bmatrix} (I - P_c \tilde{X}C),$$

and it follows from $Q_o P_c = 0$ and $Q_c P_c = I$ that

$$A(I - KC) = \begin{bmatrix} P_o & P_c \end{bmatrix} \begin{bmatrix} \Lambda_1 & 0 \\ 0 & \Lambda_2 \end{bmatrix} \begin{bmatrix} Q_o \\ Q_c - \tilde{X}C \end{bmatrix}.$$

Accordingly, the characteristic polynomial of matrix A(I - KC), denoted by det $(\lambda I - A(I - KC))$, can be given as follows:

$$\det (\lambda I - A(I - KC))$$

$$= \det \left(\left[P_o \mid P_c \right] \lambda I \left[\frac{Q_o}{Q_c} \right] - \left[P_o \mid P_c \right] \left[\frac{\Lambda_1 \mid 0}{0 \mid \Lambda_2} \right] \left[\frac{Q_o}{Q_c - \tilde{X}C} \right] \right)$$

$$= \det \left(\left[P_o \mid P_c \right] \left[\frac{(\lambda I - \Lambda_1)Q_o}{(\lambda I - \Lambda_2)Q_c + \Lambda_2 \tilde{X}C} \right] \right)$$

$$= \det(P) \det \left(\left[\frac{(\lambda I - \Lambda_1)Q_o}{(\lambda I - \Lambda_2)Q_c + \Lambda_2 \tilde{X}C} \right] \right).$$

Algorithm 1 The algorithm for generating FDIAs

1: Initialize:

Decompose matrix A in (1) as the Jordan normal form (10), Choose arbitrarily a scalar

 $\sigma \in (0,1)$ and the positive scalar M: 2: Determine the integers t, r and q according to Lemma 3, (17) and (18), respectively; 3: Set $\bar{t}_r(0) = 0$; 4: while k > 0 do 5: if $\operatorname{Re}\{\lambda_q \bar{t}_r(k)\} \ge 0$ then Set $\sigma(k+1) = \sigma$; 6: Set the attack $a(k+1) = CA\Delta \hat{x}(k) + \sigma(k+1)MI_m^t$; 7: else 8: 9: Set $\sigma(k+1) = -\sigma$; 10: Set the attack $a(k+1) = CA\Delta \hat{x}(k) + \sigma(k+1)MI_m^t$; end if 11: Calculate $\Delta \hat{x}(k+1)$ according to (7); 12: 13: Calculate $\bar{t}_r(k+1)$ according to (20); 14: k = k + 1: 15: end while

Setting $\lambda = \lambda_i$ $(i \in \{1, ..., l\})$, we can see that the last row of matrix $\lambda I - J_i$ is a zero row, which implies that there is at least a zero row in the sub-matrix $(\lambda I - \Lambda_1)Q_o$ and hence det $(\lambda I - A(I - KC)) = 0$. In other words, we conclude that λ_i (i = 1, ..., l) is the eigenvalue of matrix A(I - KC). Noting that $|\lambda_i| \ge 1$ (i = 1, ..., l), this conclusion contradicts the inequality $\rho(A(I - KC)) < 1$. As a result, there is no solution to the matrix equation (12) and the proof is complete.

From Lemma 2, the following lemma can be easily obtained.

Lemma 3: For the system (1) with estimator (2), let $\rho(A) \ge 1$ and $E_{s,t}$ represent the element of matrix E in the sth row and tth column. Define matrix $E = P^{-1}K$. Then, there exists at least one non-zero component in matrix E, that is, there exist integers $s \in \{1, \ldots, \bar{l}\}$ and $t \in \{1, \ldots, m\}$ with $\bar{l} \triangleq \sum_{i=1}^{l} n_i$ such that $E_{s,t} \neq 0$.

Proof: Let us prove the lemma by contradiction. Assume that $E_{s,t} = 0$, $\forall s \in \{1, \dots, \bar{l}\}$, $\forall t \in \{1, \dots, m\}$. That is, $E = \begin{bmatrix} 0\\ \bar{E} \end{bmatrix}$ where $\bar{E} \in \mathbb{C}^{(n-\bar{l})\times m}$ is the sub-matrix forming by the last $n-\bar{l}$ rows of E. Then, the equation K = PE can be rewritten as follows:

$$K = PE = \left[P_o \middle| P_c \right] \left[\frac{0}{\bar{E}} \right] = P_c \bar{E}.$$

The above equation implies that \overline{E} is the solution of equation (12), which contradicts the statement in Lemma 2 that equation (12) has no solution. The proof is now complete.

Before we present the necessary and sufficient condition under which the system (1) with estimator (2) is *insecure*, a procedure for generating a certain sequence of FDIAs is outlined in Algorithm 1.

Theorem 1: Assume that the attacker is able to attack all communication channels, that is, $B_a = I_m$. The system (1) with state estimator (2) is *insecure* if and only if $\rho(A) \ge 1$.

Proof: (Sufficiency) We start by proving that, if $\rho(A) \ge 1$, the system (1) state estimator (2) is *insecure*. According to Definition 1, we need to prove that there exists at least one FDIA sequence satisfying both (8) and (9) if $\rho(A) \ge 1$. In the following, we prove that (8) and (9) are true under the attacks generated by Algorithm 1.

According to Algorithm 1, it is known that

$$a(k+1) = CA\Delta\hat{x}(k) + \sigma(k+1)MI_m^t$$
(13)

PREPRINT

where $\sigma(k+1)$ takes value on either σ or $-\sigma$ with $\sigma \in (0,1)$. It follows from (6) and (13) that

$$\Delta z(k+1) = \sigma(k+1)MI_m^t,\tag{14}$$

from which we can easily see that $\|\Delta z(k+1)\| = \sigma M < M$, and this implies that condition (9) is satisfied.

To show that the condition (8) is satisfied, we define vector $t(k) = Q\Delta \hat{x}(k)$ where $t(k) = [t_1^T(k), \dots, t_p^T(k)]^T$ with $t_i(k) \in \mathbb{C}^{n_i}$ $(i \in \{1, 2, \dots, p\})$. Based on (7), (11) and Lemma 3, the dynamics of t(k) can be derived as follows:

$$t(k+1) = Jt(k) + QK\Delta z(k+1) = Jt(k) + E\Delta z(k+1).$$
(15)

Substituting (14) into (15) gives

$$t(k+1) = Jt(k) + \sigma(k+1)MEI_m^t.$$

Define $\bar{t}(k) = \left[t_1^T(k), \dots, t_l^T(k)\right]^T$ and $\underline{t}(k) = \left[t_{l+1}^T(k), \dots, t_p^T(k)\right]^T$. Noting that $J = \left[\frac{\Lambda_1 \mid 0}{0 \mid \Lambda_2}\right]$, one has
 $\bar{t}(k+1) = \Lambda_1 \bar{t}(k) + \sigma(k+1)Md,$ (16)

where $d = \left[I_{\bar{l}}, 0_{\bar{l} \times (n-\bar{l})}\right] EI_m^t$, *i.e.*, vector d is formed by the first \bar{l} elements of the tth column of matrix E. From Lemma 3, it is known that $d \neq 0$.

Define $d = \begin{bmatrix} d_1, \dots, d_{\bar{l}} \end{bmatrix}^T$ and

$$r = \underset{1 \le j \le \bar{l}}{\operatorname{argmax}} \ (d_j \ne 0), \tag{17}$$

that is, d_r is the non-zero element of vector d with the maximal index. Since $1 \le r \le \overline{l}$, there exists an integer q $(1 \le q \le l)$ such that

$$\sum_{i=1}^{q} n_i - n_q < r \le \sum_{i=1}^{q} n_i.$$
(18)

It follows from (16) that

$$\begin{bmatrix} \bar{t}_r(k+1) \\ \bar{t}_{r+1}(k+1) \\ \vdots \\ \bar{t}_{n_q}(k+1) \end{bmatrix} = \begin{bmatrix} \lambda_q & 1 \\ \lambda_q & \ddots \\ & \ddots & 1 \\ & & \lambda_q \end{bmatrix} \begin{bmatrix} \bar{t}_r(k) \\ \bar{t}_{r+1}(k) \\ \vdots \\ \bar{t}_{n_q}(k) \end{bmatrix} + \sigma(k+1)M \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix},$$
(19)

where $\bar{t}_j(k)$ is the *j*th element of vector $\bar{t}(k)$, $j \in \{r, r+1, \ldots, n_q\}$.

Noting the initial condition $\bar{t}_{r+1}(0) = 0$, it can be easily derived from (19) that $\bar{t}_{i+1}(k) = 0$ and

$$\bar{t}_r(k+1) = \lambda_q \bar{t}_r(k) + \sigma(k+1)M,$$
(20)

and therefore

$$|\bar{t}_r(k+1)|^2 = |\lambda_q|^2 |\bar{t}_r(k)|^2 + \sigma^2(k+1)M^2 + 2\sigma(k+1)M\operatorname{Re}\{\lambda_q \bar{t}_r(k)\}$$
(21)

According to Algorithm 1, it is known that $\sigma(k+1)\operatorname{Re}\{\lambda_q \bar{t}_i(k)\} \ge 0$ and $\sigma^2(k+1) = \sigma^2$. Furthermore, noticing that $|\lambda_q| \ge 1$, we have

$$|\bar{t}_r(k+1)|^2 \ge |\lambda_q|^2 |\bar{t}_r(k)|^2 + \sigma^2 M^2 \ge |\bar{t}_r(k)|^2 + \sigma^2 M^2.$$
(22)

Based on the inequality $|\bar{t}_r(k+1)|^2 \ge |\bar{t}_r(k)|^2 + \sigma^2 M^2$ and the initial condition $\bar{t}_r(0) = 0$, it can be inferred that $|\bar{t}_r(k+1)|^2 \ge (k+1)\sigma^2 M^2$, which implies that $\lim_{k\to\infty} |\bar{t}_r(k+1)| = \infty$ and therefore $\lim_{k\to\infty} t(k+1) = \infty$. Since $t(k+1) = Q\Delta \hat{x}(k+1)$, it can be deduced that at least one component of vector $\Delta \hat{x}(k+1)$ is unbounded,

and $\lim_{k\to\infty} \|\Delta \hat{x}(k+1)\| = \infty$. To this end, the condition (8) is satisfied and we finally reach the conclusion that the system is *insecure* under the attacks generated by Algorithm 1 if $\rho(A) \ge 1$.

(Necessity). To prove the necessity, we just need to show that the system \mathcal{P} in (1) with estimator \mathcal{E} in (2) is *secure* if matrix $\rho(A) < 1$. Again, let us prove by contradiction. Assume that the system (1) with estimator (2) is *insecure*, that is, there exist attacks sequences satisfying (8) and (9). It follows from (9) that $\Delta z(k+1)$ is norm-bounded. Since $\rho(A) < 1$, based on the equation $\Delta \hat{x}(k+1) = A\Delta \hat{x}(k) + K\Delta z(k+1)$, it can be inferred that $\Delta \hat{x}(k+1)$ is norm-bounded as well. That is, condition (8) is violated and the proof is now complete.

Remark 1: In the main results of [11], [14], it has been stated that the necessary and sufficient conditions for the state error by FDIAs to be unbounded are that a) the system matrix A should be unstable; and b) at least one eigenvector v corresponding to the unstable system mode satisfies $v \in Q_{oa}$ where Q_{oa} is the controllability matrix associated with the pair $(A - KCA, KB_a)$. Note that when $B_a = I_m$, condition b) is actually unnecessary and has been removed in Theorem 1 of this paper.

B. Case 2: $B_a \neq I_m$

In this case, we assume that the attacker is able to inject false data to only a part of (rather than all) communication channels, i.e., $Rk\{B_a\} < m$. It can be easily seen from Theorem 1 that, if $\rho(A) < 1$, the system (1) with estimator (2) is *secure* no matter how many communication channels the attacker could hijack. As such, in this subsection, we only consider the case when $\rho(A) \ge 1$.

The following lemmas are useful in subsequent analysis.

Lemma 4: [3] Let $A \in \mathbb{C}^{n \times m}$, $B \in \mathbb{C}^{m \times l}$ and $C \in \mathbb{C}^{k \times n}$. Assume that B has full row rank and C has full column rank. Then, $Rk\{AB\} = Rk\{A\} = Rk\{CA\}$.

Lemma 5: If the system \mathcal{P} in (1) with estimator \mathcal{E} in (2) is *insecure*, then 1) the attack sequence $\{a_k\}$ leading to the insecurity is unbounded, and 2) the state estimation difference $\Delta \hat{x}(k)$ can be represented in the following form:

$$\Delta \hat{x}(k) = P_o \zeta_1(k) + P_c \zeta_2(k) \tag{23}$$

for some $\zeta_1(k) \in \mathbb{C}^{\bar{l}}$ satisfying $\lim_{k\to\infty} \zeta_1(k) = \infty$ and some bounded vector sequence $\zeta_2(k) \in \mathbb{C}^{n-\bar{l}}$, where P_o and P_c are defined in (11).

Proof: Assume that the attack sequence $\{a_k\}$ leading to the insecurity is bounded. Noting that $\rho((I-KC)A) < 1$, it follows from the dynamics of $\Delta \hat{x}(k)$ in (7) that $\Delta \hat{x}(k+1)$ is bounded. According to Definition 1, the boundedness of $\Delta \hat{x}(k+1)$ contradicts the insecurity assumption of this lemma. As such, the attack sequence $\{a_k\}$ is unbounded.

Next, we proceed to prove that $\Delta \hat{x}(k)$ can be represented as (23) and we use the same notations for P,Q, P_o , P_c , Q_o and Q_c as defined in (10)-(11). Similar to the proof of Theorem 1, we define vector $t(k) \triangleq Q\Delta \hat{x}(k)$ and write $t(k) = \left[t_1^T(k), \ldots, t_p^T(k)\right]^T$ with $t_i(k) \in \mathbb{C}^{n_i}$ $(i \in \{1, 2, \ldots, p\})$. According to (11), the dynamics of t(k) can be given by

$$t(k+1) = \left[\frac{\bar{t}(k+1)}{\underline{t}(k+1)}\right] = \left[\frac{\Lambda_1}{0} \frac{0}{\Lambda_2}\right] \left[\frac{\bar{t}(k)}{\underline{t}(k)}\right] + \left[\frac{Q_o K}{Q_c K}\right] \Delta z(k+1),$$
(24)

where $\bar{t}(k) \triangleq \left[t_1^T(k), \dots, t_l^T(k)\right]^T$ and $\underline{t}(k) \triangleq \left[t_{l+1}^T(k), \dots, t_p^T(k)\right]^T$. As $\rho(\Lambda_2) < 1$ and $\Delta z(k)$ is norm-bounded, it can be inferred that $\underline{t}(k)$ is norm-bounded. On the other hand,

it is easy to see that $\Delta \hat{x}(k) = Pt(k) = \left[P_o \mid P_c\right] \left[\frac{\bar{t}(k)}{\underline{t}(k)}\right] = P_o \bar{t}(k) + P_c \underline{t}(k)$. Since $P_c \underline{t}(k)$ is bounded and

 $\lim_{k\to\infty} \Delta \hat{x}(k) = \infty$, it follows that $\lim_{k\to\infty} \bar{t}(k) = \infty$ and therefore expression (23) holds for $\zeta_1(k) = \bar{t}(k)$ and $\zeta_2(k) = t(k)$, which completes the proof.

Theorem 2: For the system \mathcal{P} in (1), assume that $\rho(A) \ge 1$, $\text{Rk}\{CP_o\} = s$ and the attacker is able to inject FDIAs to a part of (but not all) communication channels, *i.e.*, $\text{Rk}\{B_a\} < m$, where P_o is defined in (11). The system \mathcal{P} in (1) with estimator \mathcal{E} in (2) is secure if the following condition holds:

$$\operatorname{Rk}\left\{(I - B_a)CP_o\right\} = s. \tag{25}$$

Proof: Again, we prove the theorem by contradiction. Suppose that the system is *insecure* when condition (25) holds. It follows from Lemma 5 that (23) is true. Furthermore, noting that $\Delta z(k+1)$ is bounded, it follows from (6) and (23) that

$$a(k+1) = CA\Delta \hat{x}(k) + \Delta z(k+1) = CP_o\Lambda_1\zeta_1(k) + O(k),$$
(26)

where $O(k) \triangleq CP_c \Lambda_2 \zeta_2(k) + \Delta z(k+1)$ which is bounded.

Define matrix $\Phi = \left[\phi_1, \ldots, \phi_{\bar{l}}\right] = CP_o$ where the vector ϕ_i is equal to the *i*th column of the matrix CP_o $(1 \le i \le \bar{l})$. Since $\operatorname{Rk}\{CP_o\} = s$, there exists a matrix $\Psi = \left[\phi_{i_1}, \phi_{i_2}, \ldots, \phi_{i_s}\right]$ satisfying $\operatorname{Rk}\{\Psi\} = s$ where $1 \le i_1 < i_2 \le \ldots < i_s \le \bar{l}$. Moreover, the matrix CP_o can be represented as $CP_o = \Psi X$ where $X \in \mathbb{C}^{s \times \bar{l}}$. It can be easily found that $\operatorname{Rk}\{X\} = s$, *i.e.*, matrix X has full row rank. As a result, (26) can be represented as follows

$$a(k+1) = \Psi \xi(k) + O(k), \tag{27}$$

where $\xi(k) = X \Lambda_1 \zeta_1(k)$.

According to Lemma 5, the attack sequence $\{a(k)\}$ is unbounded, the sequence $\{O(k)\}$ is bounded, and therefore the vector sequence $\{\xi(k)\}$ is unbounded.

Left-multiplying both sides of (27) by $I - B_a$ gives rise to

$$(I - B_a)a(k+1) = (I - B_a)\Psi\xi(k) + (I - B_a)O(k),$$

and then it follows from $a(k+1) = B_a a^0(k+1)$ and $(I - B_a)B_a = 0$ that

$$(I - B_a)\Psi\xi(k) + (I - B_a)O(k) = 0.$$
(28)

Since $(I - B_a)CP_o = (I - B_a)\Psi X$ and X is full row rank, it is known from Lemma 4 that $\text{Rk}\{(I - B_a)\Psi\} = \text{Rk}\{(I - B_a)\Psi X\} = \text{Rk}\{(I - B_a)CP_o\}$. Note the fact $\text{Rk}\{(I - B_a)\Psi\} = s$ in (25) or, in other words, the matrix $(I - B_a)\Psi$ has full column rank. As $\lim_{k\to\infty} \xi(k) = \infty$, we have $\lim_{k\to\infty} (I - B_a)\Psi\xi(k) = \infty$ that contradicts (28), and the proof is now complete.

It is known From Theorem 1 that the system \mathcal{P} in (1) with estimator \mathcal{E} in (2) is *insecure* when $\rho(A) \ge 1$. In this case, it is important to ensure the security by protecting some communication channels. The following corollary provides an efficient method on which communication channels need to protected.

Corollary 1: For the system (1), assume that $\rho(A) \ge 1$ and $\text{Rk}\{CP_o\} = s$. The system \mathcal{P} in (1) with estimator \mathcal{E} in (2) is secure if

1) s communication channels are protected;

2) $\operatorname{Rk}\left\{\left[\varphi_{i_1}^T, \cdots, \varphi_{i_s}^T\right]^T\right\} = s$, where i_1, \ldots, i_s are the indexes of the protected communication channels and φ_j is the *j*th row of matrix CP_o ($i_1 \leq j \leq i_s$).

Proof: Since the communication channels i_1, \ldots, i_s are protected (i.e., free from cyber-attacks), according to the definition of matrix B_a , it is known that $\gamma_{i_1} = \ldots = \gamma_{i_s} = 0$ and $(I - B_a)CP_o = \left[\gamma_1 \varphi_1^T, \ldots, \gamma_m \varphi_m^T\right]^T$.



(a) The estimation difference.



Fig. 2. The estimation and residual differences under FDIAs

On one hand, $\operatorname{Rk}\left\{\left[\varphi_{i_1}^T, \cdots, \varphi_{i_s}^T\right]^T\right\} = s$ implies that $\operatorname{Rk}\left\{(I - B_a)CP_o\right\} \ge s$. On the other hand, we have $\operatorname{Rk}\left\{(I - B_a)CP_o\right\} \le \operatorname{Rk}\left\{CP_o\right\} = s$. As a result, $\operatorname{Rk}\left\{(I - B_a)CP_o\right\} = s$ and it follows from Theorem 2 that the system (1) with state estimator (2) is *secure*, which completes the proof.

Remark 2: It is clear that $\text{Rk}\{CP_o\} = s \leq \overline{l}$ and it can be found from (11) that \overline{l} is the number of unstable eigenvalues of matrix A (counted up to multiplicity). As such, Corollary 1 implies that the number of communication channels that should be protected is not more than the number of unstable eigenvalues of matrix A (counted up to multiplicity).

IV. SIMULATION RESULTS

In this section, we consider the state estimation system of a flight vehicle. The system consists of a moving vehicle installed with three sensors and a remote estimator. Our purpose is to 1) analyse the security of the system, and 2) protect the system from cyber-attacks if it is insecure. The linearised discrete-time model of a simplified longitudinal flight control system is described by (1) with the following system parameters:

$$A = \begin{bmatrix} 0.9944 & -0.1203 & -0.4302 \\ 0.0017 & 0.9902 & -0.0747 \\ 0 & 0.8187 & 0 \end{bmatrix}, B = \begin{bmatrix} 0.4252 \\ -0.0082 \\ 0.1813 \end{bmatrix}, C = I_3,$$

and both the system and measurement noises are assumed to be uncorrelated zero-mean white noises with covariance $W = \text{diag}\{0.1^2, 0.1^2, 0.01^2\}$ and $R = 0.1I_3$, respectively. A stationary Kalman filter is employed in the remote estimator and a χ^2 fault detector is employed as well.

It can be computed that the eigenvalues of system matrix are 1, 0.9177, and 0.0669. According to Theorem 1, the estimation system of the flight vehicle is insecure. To confirm this conclusion via simulation, a specific deceptive FDIA sequence is generated according to Algorithm 1 where the parameters are chosen as $\sigma = 0.1, M = 2$.

Fig. 2 depicts the state estimation difference $\Delta \hat{x}(k)$ and the estimation residual difference $\Delta \hat{z}(k)$ under the designed attack sequences $\{a(k)\}$. From fig. 2, it can be seen that the sequence $\{\Delta \hat{x}(k)\}$ diverges to ∞ while the sequence $\{\|\Delta \hat{z}(k)\|\}$ is always less the prescribed scalar M. Here, the estimated trajectory of the vehicle under the designed FDIA attacks deviates significantly from its nominal one but this cannot be detected by the χ^2 fault detector.

Next, let us consider how to protect the system from cyber-attacks. It can be computed that the eigenvector corresponding to the unstable system eigenvalue 1 is $P_o = \begin{bmatrix} 50.9530 & 1.2214 & 1.0000 \end{bmatrix}^T$. Since $Rk(CP_o) = 1$,

according to Corollary 1, it is known that the state estimate system of the flight vehicle is secure if the communication channel between sensor 1 and the estimator is protected.

V. CONCLUSION

In this paper, we have considered the security issues in state estimation of networked control systems, where the adversary can inject false data into the communication channels between sensors and a remote estimator. For the case that the adversary can compromise all communication channels, a necessary and sufficient condition has been derived under which the estimation error caused by the attacks is unbounded all the time. For the case that the adversary can only compromise a part of the communication channels, a sufficient condition ensuring the security is derived as well. Moreover, a criterion on protecting a sufficient number of channels such that the estimation error is kept bounded under FDIAs has been proposed. A simulation example has been proposed to demonstrate the usefulness of the developed results and algorithms. The concept of estimation residual difference has been used as an alternative measure for residual changes under cyber-attacks in [2]. One of the future topics for our research would be the analysis of cybersecurity using the new measure of KL divergence.

REFERENCES

- [1] B. D. Anderson and J. B. Moore, Optimal Filtering. Courier Dover Publications, 2005.
- [2] C.-Z. Bai, F. Pasqualetti, and V. Gupta, "Security in stochastic control systems: Fundamental limitations and performance bounds," in 2015 American Control Conference (ACC). IEEE, 2015, pp. 195–200.
- [3] D. S. Bernstein, Matrix Mathematics: Theory, Facts, and Formulas. Princeton University Press, 2009.
- [4] T. Chen, "Stuxnet, the real start of cyber warfare?[editor's note]," IEEE Network, vol. 24, no. 6, pp. 2–3, 2010.
- [5] M. S. Chong, M. Wakaiki, and J. P. Hespanha, "Observability of linear systems under adversarial attacks," in 2015 American Control Conference (ACC). IEEE, 2015, pp. 2439–2444.
- [6] H. Fawzi, P. Tabuada, and S. Diggavi, "Secure estimation and control for cyber-physical systems under adversarial attacks," *IEEE Transactions on Automatic Control*, vol. 59, no. 6, pp. 1454–146, 2014.
- [7] J. M. Hendrickx, K. H. Johansson, R. M. Jungers, H. Sandberg, and K. C. Sou, "Efficient computations of a security index for false data attacks in power networks," *IEEE Transactions on Automatic Control*, vol. 59, no. 12, pp. 3194–3208, 2014.
- [8] J. P. Hespanha, *Linear Systems Theory*. Princeton university press, 2009.
- [9] O. Kosut, "Malicious data attacks against dynamic state estimation in the presence of random noise," in *Proc. Global Conference on Signal and Information Processing*. IEEE, 2013, pp. 261–264.
- [10] O. Kosut, L. Jia, R. J. Thomas, and L. Tong, "Malicious data attacks on the smart grid," *IEEE Transactions on Smart Grid*, vol. 2, no. 4, pp. 645–658, 2011.
- [11] C. Kwon, W. Liu, and I. Hwang, "Security analysis for cyber-physical systems against stealthy deception attacks," in *Proc. American Control Conference (ACC)*. IEEE, 2013, pp. 3344–3349.
- [12] Y. Liu, P. Ning, and M. K. Reiter, "False data injection attacks against state estimation in electric power grids," in *Proc. the 16th ACM conference on Computer and Communications Security*. ACM, 2009, pp. 21–32.
- [13] S. Mishra, Y. Shoukry, N. Karamchandani, S. Diggavi, and P. Tabuada, "Secure state estimation against sensor attacks in the presence of noise," arXiv preprint arXiv:1510.02462, 2015.
- [14] Y. Mo and B. Sinopoli, "False data injection attacks in cyber physical systems," in First Workshop on Secure Control Systems, 2010.
- [15] Y. Mo, R. Chabukswar, and B. Sinopoli, "Detecting integrity attacks on scada systems," *IEEE Transactions on Control Systems Technology*, vol. 22, no. 4, pp. 1396–1407, 2014.
- [16] Y. Mo, E. Garone, A. Casavola, and B. Sinopoli, "False data injection attacks against state estimation in wireless sensor networks," in Proc. IEEE Conference on Decision and Control (CDC). IEEE, 2010, pp. 5967–5972.
- [17] Y. Mo and B. Sinopoli, "On the performance degradation of cyber-physical systems under stealthy integrity attacks," *IEEE Transactions* on Automatic Control, in press.
- [18] Y. Mo, S. Weerakkody, and B. Sinopoli, "Physical authentication of control systems: designing watermarked control inputs to detect counterfeit sensor outputs," *IEEE Control Systems*, vol. 35, no. 1, pp. 93–109, 2015.
- [19] Z.-H. Pang and G.-P. Liu, "Design and implementation of secure networked predictive control systems under deception attacks," *IEEE Transactions on Control Systems Technology*, vol. 20, no. 5, pp. 1334–1342, 2012.

- [20] F. Pasqualetti, F. Dorfler, and F. Bullo, "Attack detection and identification in cyber-physical systems," *IEEE Transactions on Automatic Control*, vol. 58, no. 11, pp. 2715–2729, 2013.
- [21] Y. Shoukry and P. Tabuada, "Event-triggered state observers for sparse sensor noise/attacks," *IEEE Transactions on Automatic Control*, vol. 61, no. 8, pp. 2079–2091, 2016.
- [22] A. Teixeira, S. Amin, H. Sandberg, K. H. Johansson, and S. S. Sastry, "Cyber security analysis of state estimators in electric power systems," in *Proc. IEEE Conference on Decision and Control (CDC)*. IEEE, 2010, pp. 5991–5998.
- [23] A. Teixeira, I. Shames, H. Sandberg, and K. H. Johansson, "A secure control framework for resource-limited adversaries," *Automatica*, vol. 51, pp. 135–148, 2015.
- [24] S. Weerakkody, X. Liu, S. H. Son, and B. Sinopoli, "A graph theoretic characterization of perfect attackability for secure design of distributed control systems," *IEEE Transactions on Control of Network Systems*, in press.
- [25] M. Zhu and S. Martínez, "On the performance analysis of resilient networked control systems under replay attacks," *IEEE Transactions on Automatic Control*, vol. 59, no. 3, pp. 804–808, 2014.