

Data-Driven Attack Detection for Linear Systems

Vishaal Krishnan and Fabio Pasqualetti

Abstract—This paper studies the attack detection problem in a data-driven and model-free setting, for deterministic systems with linear and time-invariant dynamics. Differently from existing studies that leverage knowledge of the system dynamics to derive security bounds and monitoring schemes, we focus on the case where the system dynamics, as well as the attack strategy and attack location, are unknown. We derive fundamental security limitations as a function of only the observed data and without estimating the system dynamics (in fact, no assumption is made on the identifiability of the system). In particular, (i) we derive detection limitations as a function of the informativity and length of the observed data, (ii) provide a data-driven characterization of undetectable attacks, and (iii) construct a data-driven detection monitor. Surprisingly, and in accordance with recent studies on data-driven control, our results show that model-based and data-driven security techniques share the same fundamental limitations, provided that the collected data remains sufficiently informative.

Index Terms—Data-driven security and attack detection.

I. INTRODUCTION

The increasing integration of the cyber and physical layers in many real-world systems has led to the emergence of cyber-physical systems security as a prominent engineering discipline, with attack monitoring forming a crucial component. Its methods differ from traditional information security techniques which lack an appropriate abstraction of the physical layer [1] and are inadequate for the protection of cyber-physical systems, where the target is often the underlying dynamics of the physical system.

Attack monitoring methods for cyber-physical systems can be broadly classified into model-based and data-driven approaches. While the latter relies only on the data, generated in the form of measurements from sensors deployed on the physical system, the former additionally assumes knowledge of the model of the underlying system. Clearly, from the perspective of implementation, model-based monitoring methods [2], [3] have to first contend with the difficulty of obtaining a reliable model for the underlying system. This is in practice achieved by system identification, for which the available data must be adequately informative – a requirement that is difficult to meet for complex systems. Moreover, it is unclear if full system identification is even necessary for attack monitoring. These considerations have contributed to the increasing popularity in recent years of data-driven approaches to monitoring. However, despite the proliferation of data-driven methods for security, a detailed

characterization of their limitations is lacking, especially when compared to model-based methods. This letter fills such a gap for the monitoring of linear systems.

We now provide an overview of the data-driven attack monitoring problem studied in this letter. We consider the following discrete-time linear time-invariant system:

$$\begin{aligned}x(k+1) &= Ax(k) + Bu(k), \\y(k) &= Cx(k) + Du(k),\end{aligned}\tag{1}$$

with $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$, $C \in \mathbb{R}^{p \times n}$ and $D \in \mathbb{R}^{p \times m}$ the system matrices, $x : \mathbb{N} \rightarrow \mathbb{R}^n$ the state, $u : \mathbb{N} \rightarrow \mathbb{R}^m$ the (attack) control input, and $y : \mathbb{N} \rightarrow \mathbb{R}^p$ the output of the system. Further, let $u = (u_1, \dots, u_m)$, where $u_j : \mathbb{N} \rightarrow \mathbb{R}$ is the input to actuator $j \in \{1, \dots, m\}$ and $y = (y_1, \dots, y_p)$, where $y_i : \mathbb{N} \rightarrow \mathbb{R}$ is the output from sensor $i \in \{1, \dots, p\}$.

The attacks are modeled in the form of data injection to the system (1), which includes a large class of attacks [2]. Thus, the system (1) is said to be under attack if $u \neq 0$. The data-driven attack detection problem reads as follows:

Problem 1: (Data-driven attack detection) Given the output $\{y(k)\}_{k \in \mathbb{N}}$, determine if the (attack) input $u \neq 0$. \square

An algorithm to solve Problem 1 is called *attack monitor* [2]. In Problem 1, the matrices A, B, C, D of system (1), the state $\{x(k)\}_{k \in \mathbb{N}}$, and the (attack) input $\{u(k)\}_{k \in \mathbb{N}}$ are unknown. Further, the Attack Detection Problem 1 is one of binary classification of time-series of measurements $\{y(k)\}_{k \in \mathbb{N}}$ into the classes $\{\text{Attack}, \text{No – Attack}\}$, where $u \neq 0$ corresponds to the class *Attack* and $u \equiv 0$ to *No – Attack*. Therefore, data-driven attack detection is essentially an inverse problem of determining the structure of the inputs to system (1) based on the output measurements. In this letter, we characterize the fundamental limitations of data-driven attack monitoring when the underlying system is linear and time-invariant.

Related work. Approaches to data-driven monitoring fundamentally rely on the fact that the underlying system (1) imposes structure on the temporal characteristics of the output data streams. The task of a data-driven attack monitor is then one of detecting changes in the structure of the output data streams in the absence of a model of the underlying system that generates them. Viewed this way, the task of data-driven attack monitoring is closely related to outlier detection in sensor data streams [4], [5]. Recent works have proposed machine learning methods for data-driven attack monitoring of cyber-physical systems, which broadly fall within the categories of supervised [6] and unsupervised learning [7]–[9]. While these works demonstrate performance of varying degree in implementation, a theoretical characterization of performance of these approaches remains elusive.

This material is based upon work supported in part by UCOP-LFR-18-548175 and AFOSR-FA9550-19-1-0235. Vishaal Krishnan and Fabio Pasqualetti are with the Department of Mechanical Engineering, University of California at Riverside, Riverside, CA 92521 USA. E-mail: vishaalk@ucr.edu, fabiopas@engr.ucr.edu.

The performance of data-driven methods relies heavily on the quality of the available data. Therefore, an investigation into the fundamental limitations of data-driven methods must relate notions of performance to notions of data informativity. In the context of system identification, persistency of excitation [10] is often assumed as a precondition on the data. In [11], the authors propose the notion of data informativity for data-driven control [12], [13], where they characterize the necessity of persistency of excitation for data-driven analysis and control and show that for certain problems, persistency of excitation is not necessary. In a similar vein, in this letter we obtain conditions on the information content in the measurement data in the context of data-driven monitoring.

Finally, in [2] the authors characterize undetectable attacks for model-based attack monitoring. Other works have studied the problem of false-data injection [14], [15], where the objective is to inject inputs that will remain undetected by an attack monitor. Here, we provide an equivalent characterization of undetectable attacks in a data-driven setting

Paper contributions. This letter contributes a characterization of the fundamental limitations of data-driven attack monitoring, from a systems-theoretic perspective. The particular contributions and the outline we follow are detailed below: (i) We first briefly treat the problem of attack detection in the model-based setting and develop the framework and preliminary results that we then adapt to the data-driven setting. (ii) We propose the notion of Hankel information to characterize the information content in the output time series. We then obtain a systems-theoretic bound on the information content of the output time series generated by a given system and the time taken to attain this bound. (iii) We propose a data-driven attack monitoring scheme that relies on learning the dynamics governing the features in the output data. (iv) We provide a practical data-driven heuristic for handling the output data for use in the attack monitor, and characterize the length of the time horizon for data collection to reach detection capability under the heuristic. (v) We finally characterize attacks undetectable by data-driven monitors.

II. DATA-DRIVEN ATTACK DETECTION

In this section, we address the data-driven attack detection problem. The output of system (1) over a time window $\{0, 1, \dots, N-1\}$, for any $N \in \mathbb{N}$, can be expressed as

$$\mathbf{y}_{0:N-1} = [\mathcal{O}_N \mid \mathcal{C}_N] \begin{pmatrix} x(0) \\ \mathbf{u}_{0:N-1} \end{pmatrix}, \quad (2)$$

where $\mathbf{y}_{r:s} = (y(r), y(r+1), \dots, y(s)) \in \mathbb{R}^{p(s-r+1)}$ and $\mathbf{u}_{r:s} = (u(r), u(r+1), \dots, u(s)) \in \mathbb{R}^{m(s-r+1)}$, for any $r, s \in \mathbb{N}$ such that $r \leq s$. The matrices \mathcal{O}_N and \mathcal{C}_N read as:

$$\mathcal{O}_N = \begin{bmatrix} C \\ CA \\ CA^2 \\ \vdots \\ CA^{N-1} \end{bmatrix}, \quad \mathcal{C}_N = \begin{bmatrix} D & 0 & \dots & 0 \\ CB & D & \dots & 0 \\ CAB & CB & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ CA^{N-2}B & CA^{N-3}B & \dots & D \end{bmatrix}.$$

Notice from (2) that the outputs over time-windows of size N belong to the column space of $[\mathcal{O}_N \mid \mathcal{C}_N]$.

A. Feature space dynamics and undetectable attacks

In the ensuing analysis, we consider the nominal system under no attack (i.e., $u \equiv 0$). We now factor \mathcal{O}_N by singular value decomposition [16] and construct $S^{(m, \text{nom})}$ from the left singular vectors of \mathcal{O}_N corresponding to the non-zero singular values, (where m in the superscript denotes that it is obtained from the model). Notice that $\text{Col}(S^{(m, \text{nom})}) = \text{Col}(\mathcal{O}_N)$, and let $\text{Col}(S^{(m, \text{nom})})$ be the *feature space* of the nominal outputs of (1) over time windows of size N . Since the columns of $S^{(m, \text{nom})}$ are orthonormal, they form a basis of the nominal feature space. Then, the nominal output over a time window $\{k, \dots, k+N-1\}$ is $\mathbf{y}_{k:k+N-1} = \mathcal{O}_N x(k)$, and it can be expressed in the feature space coordinates as

$$\mathbf{w}(k) = S^{(m, \text{nom})\top} \mathbf{y}_{k:k+N-1} = S^{(m, \text{nom})\top} \mathcal{O}_N x(k).$$

We call $\{\mathbf{w}(k)\}_{k \in \mathbb{N}}$ the nominal *feature vector sequence*. We note that $\{\mathbf{w}(k)\}_{k \in \mathbb{N}}$ is the representation of the nominal outputs over windows of size N of (1) in the feature space coordinates. It can be shown, for sufficiently large N , that the nominal feature vector sequence is generated by the system:

$$\mathbf{w}(k+1) = \left(S^{(m, \text{nom})\top} \mathcal{O}_N \right) A \left(S^{(m, \text{nom})\top} \mathcal{O}_N \right)^\dagger \mathbf{w}(k), \quad (3)$$

where $\mathbf{w}(0) = \mathbf{w}_0 = S^{(m, \text{nom})\top} \mathcal{O}_N x(0)$. This forms the content of the Theorem 2.1 below, which is proved in Appendix I. Before stating the theorem, we recall the notion of observability index of linear time-invariant systems, which provides a lower bound on the size N of the time window. The observability index of the system (1) is defined as:

$$\nu = \min \{ N \in \mathbb{N} \mid \text{Rank}(\mathcal{O}_N) = \text{Rank}(\mathcal{O}_{N+j}), \forall j \in \mathbb{N} \}.$$

Theorem 2.1: (Feature space dynamics) Let ν be the observability index of (1). For any $N \geq \nu$, the feature vector sequence $\{\mathbf{w}(k)\}_{k \in \mathbb{N}}$ generated by the system (1) with $u \equiv 0$ is a solution to (3). \square

Theorem 2.1 suggests the design of an attack detection scheme that is based on verifying if the output time series $\{y(k)\}_{k \in \mathbb{N}}$, generated by the system (1), can be completely characterized by a feature vector sequence $\{\mathbf{w}(k)\}_{k \in \mathbb{N}}$ that is a solution to (3). In particular, if $\mathbf{y}_{k:k+N-1} = S^{(m, \text{nom})} \mathbf{w}(k)$ for all $k \in \mathbb{N}$, where $\{\mathbf{w}(k)\}_{k \in \mathbb{N}}$ is a solution to (3), then $\{y(k)\}_{k \in \mathbb{N}}$ is classified as No-Attack. The output sequence is otherwise classified as Attack. We also note that maximum detection capability is attained at time $T_{\text{safe}}^m = \nu$.

Remark 1: (Undetectable attacks) We first note that for time instants $T \leq \nu$, an (attack) input $\mathbf{u}_{0:T-1}$ over the time horizon $\{0, \dots, T-1\}$ is detectable if and only if $\mathcal{C}_T \mathbf{u}_{0:T-1} \notin \text{Col}(\mathcal{O}_T)$. Therefore, any (attack) input $\mathbf{u}_{0:T-1}$ for which $\mathcal{C}_T \mathbf{u}_{0:T-1} \in \text{Col}(\mathcal{O}_T)$ is undetectable. However, for attacks starting at a time instant $T \geq \nu + 1$ (i.e., $u : \mathbb{N} \rightarrow \mathbb{R}^m$ such that $\mathbf{u}_{0:T-1} = 0$), not only must the outputs $\mathbf{y}_{k:k+N-1} \in \text{Col}(\mathcal{O}_N)$, they must also, for $N \geq \nu + 1$, be completely characterized by a feature vector sequence $\{\mathbf{w}(k)\}_{k \in \mathbb{N}}$ that is a solution to (3). Thus, an attack starting at $T \geq \nu + 1$ is undetectable if and only if $\mathbf{u}_{T:T+N-1} \in \text{Ker}(\mathcal{C}_N)$ for any $N \in \mathbb{N}$.

Furthermore, if $u : \mathbb{N} \rightarrow \mathbb{R}^m$ is a finite duration attack (i.e., there exists a $d \in \mathbb{N}$ such that $\mathbf{u}_{T+d:\infty} = 0$), it can remain undetectable (i.e., $\mathbf{u}_{T:T+N-1} \in \text{Ker}(\mathcal{C}_N)$ for any $N \in \mathbb{N}$) only if the system (1) is not left invertible [17].

This analysis is compatible with the fundamental limitations derived for model-based attack detection in [2] in terms of the zero dynamics of (1). \square

B. Information bound on output time series

The preceding analysis assumes the knowledge of the system matrix A and the observability matrix \mathcal{O}_N of system (1), which are unknown in the data-driven setting. We instead have access only to a finite time series of outputs over a horizon $\{0, \dots, T-1\}$, from which we can construct the following Hankel matrix with a time window of size $N \leq T$:

$$Y_{N,T} = \begin{bmatrix} y(0) & y(1) & \dots & y(T-N) \\ y(1) & y(2) & \dots & y(T-N+1) \\ y(2) & y(3) & & \\ \vdots & \vdots & \ddots & \vdots \\ y(N-1) & y(N) & \dots & y(T-1) \end{bmatrix}. \quad (4)$$

As a first step, we characterize the information bound on the output time series for the attack detection problem via an upper bound on the rank of $Y_{N,T}$. Notice that $\text{Rank}(Y_{N,T})$ is a function of $N, T \in \mathbb{N}$ and the initial state $x(0)$ (since $\{y(k)\}_{k \in \mathbb{N}}$ is generated by the underlying system (1)). We introduce the notion of Hankel information defined below:

$$\Gamma(\{y(k)\}_{k \in \mathbb{N}}) = \sup_{N, T \in \mathbb{N}} \{\text{Rank}(Y_{N,T}) : N \leq T\}, \quad (5)$$

as the measure of the information content in the output data $\{y(k)\}_{k \in \mathbb{N}}$. In fact, we have $Y_{k,N,T} = \mathcal{O}_N [x(k) \ Ax(k) \ \dots \ A^{T-N}x(k)]$, from which it follows that $\text{Rank}(Y_{N,T}) \leq \text{Rank}(\mathcal{O}_N)$. When $\text{Rank}(Y_{N,T}) = \text{Rank}(\mathcal{O}_N)$, we can completely reconstruct the column space of \mathcal{O}_N from $\{y(k)\}_{k \in \mathbb{N}}$. This corresponds to the full information scenario. However, since this is not generally the case, we provide a bound in Proposition 2.2 on the Hankel information of the output data $\{y(k)\}_{k \in \mathbb{N}}$ generated from an initial state $x(0)$. Before stating the proposition, we first introduce the notion of excitability index for the system 1.

Definition 1: (Excitability index) Let $x \in \mathbb{R}^n$ and $N \in \mathbb{N}$, and let $E_N(x) = [x \ Ax \ \dots \ A^{N-1}x]$. The excitability index $\mu(x)$ of (1) at x is defined as

$$\mu(x) = \min\{i : \text{Rank}(E_i(x)) = \text{Rank}(E_{i+j}(x)), \forall j \in \mathbb{N}\}.$$

The excitability index of (1) is $\mu = \max_{x \in \mathbb{R}^n} \mu(x)$. \square

Proposition 2.2: (Hankel information bound on output time series) Let $\{y(k)\}_{k \in \mathbb{N}}$ be output of (1) from an initial state $x(0)$. Then, for any $N, T \in \mathbb{N}$ with $N \leq T$, we have:

- (i) $\text{Rank}(Y_{N,T}) \leq \min_{\mathcal{V} \subseteq \text{Ker}^\perp(\mathcal{O}_N)} \{\dim(\mathcal{V}) : x(0) \in \mathcal{V}; A\mathcal{V} \subseteq \mathcal{V}\}$,
- (ii) $\Gamma(\{y(k)\}_{k \in \mathbb{N}}) \leq \min\{\nu p, \mu\}$,

where $Y_{N,T}$ and $\Gamma(\{y(k)\}_{k \in \mathbb{N}})$ are as defined in (4) and (5).

Proof: We start by noticing that $Y_{N,T} = \mathcal{O}_N [x(0) \ Ax(0) \ \dots \ A^{T-N}x(0)]$, from which it follows that $\text{Rank}(Y_{N,T}) \leq \text{Rank}([x(0) \ Ax(0) \ \dots \ A^{T-N}x(0)])$. Let \mathcal{V} be A -invariant, and let $x(0) \in \mathcal{V} \subseteq \text{Ker}^\perp(\mathcal{O}_N)$. Then, $\dim(\{x(0), Ax(0), \dots, A^{T-N}x(0)\}) \leq \dim(\mathcal{V})$. It then follows that $\text{Rank}(Y_{N,T}) \leq \text{Rank}([x(0) \ Ax(0) \ \dots \ A^{T-N}x(0)]) \leq \dim(\mathcal{V})$. Since the above inequality holds for any A -invariant $\mathcal{V} \subseteq \text{Ker}^\perp(\mathcal{O}_N)$, the claim follows.

Moreover, we have $\text{Rank}(\mathcal{O}_N) \leq \text{Rank}(\mathcal{O}_\nu)$ and $\text{Rank}([x(0) \ Ax(0) \ \dots \ A^{T-N}x(0)]) \leq \text{Rank}([x(0) \ Ax(0) \ \dots \ A^{\mu-1}x(0)])$ where ν and μ are the observability and excitability indices respectively. Moreover, we have $\text{Rank}(Y_{N,T}) = \dim(\text{Ker}^\perp(\mathcal{O}_N) \cap \text{Col}([x(0) \ Ax(0) \ \dots \ A^{T-N}x(0)]))$. Since $\dim(\text{Ker}^\perp(\mathcal{O}_N)) = \dim(\text{Col}(\mathcal{O}_N)) \leq \min\{\nu p, n\}$, $\dim(\text{Col}([x(0) \ Ax(0) \ \dots \ A^{T-N}x(0)])) \leq \mu$ for any $N, T \in \mathbb{N}$ and $x(0) \in \mathbb{R}^n$, and we have $\mu \leq n$, we get $\sup\{\text{Rank}(Y_{N,T}) : N, T \in \mathbb{N} \text{ and } N \leq T\} \leq \min\{\nu p, \mu\}$. \blacksquare

Proposition 2.2 implies that, if the initial state $x(0)$ belongs to an A -invariant subspace of $\text{Ker}^\perp(\mathcal{O}_N)$ of lower dimension, then the column space of \mathcal{O}_N cannot be completely characterized from the Hankel matrix $Y_{N,T}$, for any $T \in \mathbb{N}$. Proposition 2.2 therefore characterizes an information bound on $\{y(k)\}_{k \in \mathbb{N}}$. We now characterize the length of the shortest time horizon over which $\{y(k)\}_{k \in \mathbb{N}}$ attains the Hankel information bound.

Theorem 2.3: (Minimum horizon length) Let $\{y(k)\}_{k \in \mathbb{N}}$ be the outputs of (1) from the initial state $x(0)$. Then, for any $x(0) \in \mathbb{R}^n$ and $T_0 \geq \nu + \mu - 1$, there exists $N_0 \in \mathbb{N}$ with $N_0 \leq T_0$ such that:

$$\text{Rank}(Y_{N_0, T_0}) = \Gamma(\{y(k)\}_{k \in \mathbb{N}}),$$

where $Y_{N,T}$ and $\Gamma(\{y(k)\}_{k \in \mathbb{N}})$ are as defined in (4) and (5).

Proof: For the Hankel matrix $Y_{N,T}$, we have $r = Np$ rows and let c be the number of columns. Then by construction, we have $T = N + c - 1$. Thus, we have $Y_{N,T} = \mathcal{O}_N [x(0) \ Ax(0) \ \dots \ A^{c-1}x(0)]$, and $\text{Rank}(Y_{N,T}) = \dim(\text{Ker}^\perp(\mathcal{O}_N) \cap \text{Col}([x(0) \ Ax(0) \ \dots \ A^{c-1}x(0)]))$. We also have $\text{Ker}^\perp(\mathcal{O}_N) \subseteq \text{Ker}^\perp(\mathcal{O}_\nu)$ and $\text{Col}([x(0) \ Ax(0) \ \dots \ A^{c-1}x(0)]) \subseteq \text{Col}([x(0) \ Ax(0) \ \dots \ A^{\mu-1}x(0)])$ for all $N, c \in \mathbb{N}$. Therefore, we get $\text{Ker}^\perp(\mathcal{O}_N) \cap \text{Col}([x(0) \ Ax(0) \ \dots \ A^{c-1}x(0)]) \subseteq \text{Ker}^\perp(\mathcal{O}_\nu) \cap \text{Col}([x(0) \ Ax(0) \ \dots \ A^{\mu-1}x(0)])$. Thus, for $T_0 \geq \nu + \mu - 1$, we can choose $N_0 = \nu$ and $c = \mu$ such that Y_{N_0, T_0} attains maximum rank. \blacksquare

Theorem 2.3 states that for any initial state $x(0)$, we attain the Hankel information (upper) bound in Proposition 2.2 for the time series of outputs $\{y(k)\}_{k \in \mathbb{N}}$ generated by the system (1), within the finite time horizon $\{0, \dots, \nu + \mu - 1\}$.

C. Data-driven attack detection monitor

For the nominal system under no attack (i.e., $u \equiv 0$), we now seek to obtain a data-driven expression for the feature space dynamics, as in (3). We factor the Hankel matrix $Y_{N,T} \in \mathbb{R}^{pN \times (T-N+1)}$ by singular value decomposition to obtain a matrix $S_{N,T} \in \mathbb{R}^{pN \times q}$ whose column vectors are the (orthonormal) left singular vectors of $Y_{N,T}$ corresponding to the non-zero singular values. The columns of $S_{N,T}$ can be interpreted as the *features in the output data* $\{y(k)\}_{k=0}^{T-1}$ with a time window of size N . In the ensuing analysis, we fix N, T and suppress them from the notation, hereby denoting $Y_{N,T}$ and $S_{N,T}$ simply by Y and S , respectively.

We now define feature vectors $\mathbf{w}(k) = S^\top \mathbf{y}_{k:k+N-1}$ by projecting the outputs $\mathbf{y}_{k:k+N-1}$, for $k \in \{0, \dots, T-N\}$, onto the feature space, and construct a matrix W as follows:

$$W = \begin{bmatrix} | & | & \dots & | \\ \mathbf{w}(0) & \mathbf{w}(1) & \dots & \mathbf{w}(T-N) \\ | & | & \dots & | \end{bmatrix}.$$

We then similarly construct a matrix \vec{W} from output data over a time horizon $\{1, \dots, T\}$, i.e., $\{y(k)\}_{k=1}^T$, to get:

$$\vec{W} = \begin{bmatrix} | & | & \dots & | \\ \mathbf{w}(1) & \mathbf{w}(2) & \dots & \mathbf{w}(T-N+1) \\ | & | & \dots & | \end{bmatrix}.$$

We show in Theorem 2.4 that, for T greater than the minimum horizon length in Theorem 2.3, the nominal feature vector sequence is obtained as a solution to $\mathbf{w}(k+1) = M^* \mathbf{w}(k)$ for all $k \in \mathbb{N}$, where $M^* \in \mathbb{R}^{q \times q}$ is a solution to the following minimization problem:

$$M^* = \arg \min_{M \in \mathbb{R}^{q \times q}} \|\vec{W} - MW\|_F, \quad (6)$$

where $\|\cdot\|_F$ denotes the Frobenius norm. In the data-driven analysis literature, this approach goes by the name of Dynamic Mode Decomposition [18], and is closely related to the Eigensystem Realization Algorithm [19] in the domain of system identification. Theorem 2.4 characterizes the conditions under which the feature vector sequence offers a complete characterization of the nominal output time series.

Theorem 2.4: (Data-driven feature space dynamics) For $N \geq \nu$ and $T \geq N + \mu - 1$, the feature vector sequence $\{\mathbf{w}(k)\}_{k \in \mathbb{N}}$ generated by the (nominal) system (1) with $u \equiv 0$ is a solution to $\mathbf{w}(k+1) = M^* \mathbf{w}(k)$ for all $k \in \mathbb{N}$, where M^* is given by (6). Moreover, the global minimizer in (6) is $M^* = \vec{W}W^\dagger$, and the minimum value is 0.

The proof of Theorem 2.4 is provided in Appendix II. It follows from Theorem 2.4 that, if the system is not under attack over the time horizon $\{0, 1, \dots, T\}$ with $T \geq \nu + \mu - 1$, the feature vector sequence $\{\mathbf{w}(k)\}_{k=0}^T$ is a solution to $\mathbf{w}(k+1) = M^* \mathbf{w}(k)$, where M^* can be computed from the (nominal) output sequence $\{y(k)\}_{k=0}^T$ via the minimization (6). This enables the prediction of the outputs starting at time $T+1$, as $\hat{\mathbf{y}}_{k:k+N-1} = S\mathbf{w}(k)$ for $k \in \{T-N+1, \dots\}$. An attack detection scheme can therefore be defined based on the prediction error, where the case $\mathbf{y}_{k:k+N-1} \neq \hat{\mathbf{y}}_{k:k+N-1}$ is classified as Attack. Moreover, the maximum data-driven

attack detection capability is attained at time instant $T_{\text{safe}} = \nu + \mu$. In comparison, in the model-based setting attack detection capability is attained at time $T_{\text{safe}}^m = \nu$, where ν is the observability index, as established in Theorem 2.1 and discussed in Remark 1. We refer the reader to [2] for detailed accounts on model-based attack monitor design.

Although attack detection capability can be attained at $T_{\text{safe}} = \nu + \mu$, it requires the construction of a Hankel matrix of appropriate dimensions (with νp rows and μ columns). However, since ν and μ are unknown in the data-driven setting, achieving full detection capability at $T_{\text{safe}} = \nu + \mu$ requires an algorithm to compute N , which is likely to increase the time complexity of the monitor. We therefore provide a practical data-driven heuristic in Remark 2 for the choice of N , and an estimate in Theorem 2.5 for the time \bar{T}_{safe} at which full attack detection capability is attained for the heuristic.

Remark 2: (Data-driven choice of N) Given a finite time series of outputs over a time horizon $\{0, \dots, T-1\}$, and a window of size N , the Hankel matrix satisfies $Y_{N,T} \in \mathbb{R}^{pN \times (T-N+1)}$. In our heuristic algorithm, we select N to satisfy $T-N+1 \geq pN$, to maintain at least as many columns as there are rows. This leads to the choice $N = \left\lfloor \frac{T+1}{p+1} \right\rfloor$. \square

Theorem 2.5: (Safe time horizon length) Let $\{y(k)\}_{k \in \mathbb{N}}$ be the output of (1), let $N(T) = \left\lfloor \frac{T+1}{p+1} \right\rfloor$, and let $\bar{T}_{\text{safe}} \geq \max \left\{ \nu(p+1) - 1, \mu \left(\frac{p+1}{p} \right) - 1 \right\}$. Then,

$$\text{Rank} \left(Y_{N(\bar{T}_{\text{safe}}), \bar{T}_{\text{safe}}} \right) = \Gamma \left(\{y(k)\}_{k \in \mathbb{N}} \right),$$

where $Y_{N,T}$ and $\Gamma(\{y(k)\}_{k \in \mathbb{N}})$ are as defined in (4) and (5).

Proof: Following the arguments in the proof of Theorem 2.3, the observability index ν yields a necessary and sufficient lower bound on N , and we get a lower bound on T by allowing $\left\lfloor \frac{T+1}{p+1} \right\rfloor \geq \nu$. Moreover, the excitability index yields a necessary and sufficient lower bound on the number of columns as $T - N + 1 \geq \mu$. Therefore, for a safe time horizon length \bar{T}_{safe} satisfying $\bar{T}_{\text{safe}} \geq \max \left\{ \nu(p+1) - 1, \mu \left(\frac{p+1}{p} \right) - 1 \right\}$, we get that $\text{Rank} \left(Y_{N(\bar{T}_{\text{safe}}), \bar{T}_{\text{safe}}} \right)$ attains the maximum value. \blacksquare

We next characterize undetectable attacks for the data-driven attack monitor. Notice that any attack over the time horizon $\{0, \dots, \nu + \mu\}$ is undetectable, since attack detection capability is only attained at $T_{\text{safe}} = \nu + \mu$. We therefore restrict our attention to attacks starting at $T \geq T_{\text{safe}} + 1$.

Theorem 2.6: (Data-driven undetectable attacks) The (attack) input $u : \mathbb{N} \rightarrow \mathbb{R}^m$ to the system (1), with $\mathbf{u}_{0:T-1} = 0$ for $T \geq T_{\text{safe}} + 1 = \nu + \mu + 1$, is undetectable if and only if $\mathbf{u}_{T:T+N-1} \in \text{Ker}(\mathcal{C}_N)$ for all $N \in \mathbb{N}$.

Proof: We have $\mathbf{w}(T-N) = S^\top \mathbf{y}_{T-N:T-1}$ and from the feature space dynamics, we get $\mathbf{w}(T-N+1) = M^* \mathbf{w}(T-N)$, where M^* is obtained from (6) over the time horizon $\{0, \dots, T_{\text{safe}}\}$ and the predicted output over $\{T-N+1, \dots, T\}$ is given by $\hat{\mathbf{y}}_{T-N+1:T} = S\mathbf{w}(T-N+1) = SM^*S^\top \mathbf{y}_{T-N:T-1}$. It follows from (2) and Theorem 2.4 that the output prediction error $\mathbf{y}_{T-N+1:T} - \hat{\mathbf{y}}_{T-N+1:T} =$

$\mathcal{C}_N \mathbf{u}_{T-N+1:T}$, where $\mathbf{u}_{T-N+1:T} = (0, \dots, 0, u(T))$. For an undetectable attack, the prediction error $\mathbf{y}_{T-N+1:T} - \hat{\mathbf{y}}_{T-N+1:T} = \mathcal{C}_N \mathbf{u}_{T-N+1:T} = 0$, which implies that $\mathbf{u}_{T-N+1:T} \in \text{Ker}(\mathcal{C}_N)$, i.e., $u(T) \in \text{Ker}(D)$. By induction, we get that $\mathbf{u}_{T:T+N-1} \in \text{Ker}(\mathcal{C}_N)$ for any $N \in \mathbb{N}$. Conversely, if $\mathbf{u}_{T:T+N-1} \in \text{Ker}(\mathcal{C}_N)$ for any $N \in \mathbb{N}$, we see that the prediction error vanishes, which implies that the attack is undetectable, thereby proving the claim. ■

III. NUMERICAL EXPERIMENTS

We now present results from numerical experiments validating the key theoretical results presented in this letter and comparing model-based and data-driven attack monitoring. We considered a linear time-invariant system (1) of state-space dimension $n = 50$ defined below:

$$A = \left[\begin{array}{c|c} \mathbf{0}_{(n-1) \times 1} & I_{n-1} \\ \hline -1 & -\mathbf{1}_{1 \times (n-1)} \end{array} \right],$$

with $m = 5$ actuators and $p = 10$ sensors. The columns of B and the rows of C were distinctly chosen at random from $\{e_i\}_{i=1}^n$, the standard basis for \mathbb{R}^n . The matrix D was chosen to be the zero matrix.

We determined the values of the observability and excitability indices numerically to be $\nu = 15$ and $\mu = 50$, respectively. With the data-driven heuristic from Remark 2 we plot in Figure 1 the rank of the Hankel matrix $Y_{N(T),T}$, $\text{Rank}(Y_{N(T),T})$ as a function of T . We then designed

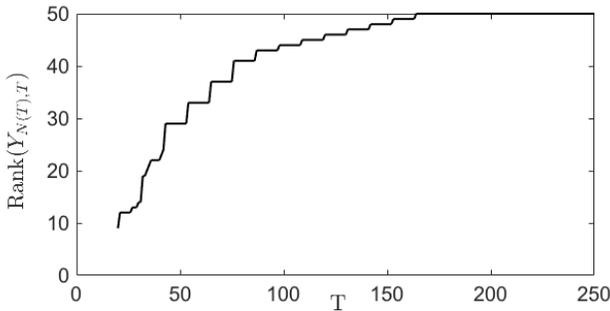


Fig. 1. The figure shows the plot of the rank of the Hankel matrix $Y_{N(T),T}$ vs. time T with N chosen according to the data-driven heuristic from Remark 2. We observe that $\text{Rank}(Y_{N(T),T})$ is a monotonically increasing function of T , and a maximum rank of 50 is attained at $\bar{T}_{\text{safe}} = 164$. With $p = 10$ and $\nu = 15$ obtained numerically, we see that $\nu(p+1) - 1 = 164 = \bar{T}_{\text{safe}}$, which validates Theorem 2.5.

model-based and data-driven attack detection monitors based on feature space dynamics (3) and as outlined in Section II-C, respectively. We injected the system with an attack input at the actuator $j = 4$ at $T = 249 > \bar{T}_{\text{safe}}$, such that $u_4(249) = 1$, and tracked the prediction error for the model-based and data-driven attack detection monitors starting from a random initial state, as shown in Figure 2.

IV. CONCLUSION AND FUTURE WORK

In this letter, we characterized the fundamental limitations on data-driven monitoring of linear time-invariant systems from a systems-theoretic perspective. In particular: (i) we

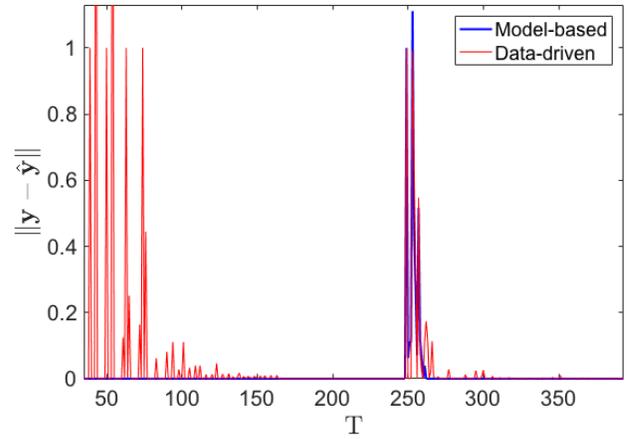


Fig. 2. The figure shows the responses of the model-based and data-driven attack monitors for an attack input $u_4(249) = 1$. Attack detection capability for the model-based monitor is achieved at time $T = 15$ and for data-driven monitor at time $\bar{T}_{\text{safe}} = 164$ (consistent with the result in Figure 1) as seen from the convergence to zero of the prediction error, implicitly validating Theorems 2.1 and 2.4. The attack is detected as an error in the output prediction by the monitors at $T = 249$. We also observe that the data-driven monitor recovers from the attack at around $T = 270$ when the prediction error is below 0.03 units ($< 3\%$ of attack input magnitude) and gradually decays to zero.

characterized an information bound on the output time series and the minimum time horizon length for measurement data collection to achieve detection capability; (ii) we provided a heuristic for the choice of dimensions of the Hankel matrix employed in data-driven attack detection; and (iii) we obtained a characterization of undetectable attacks. Surprisingly, our results show that model-based and data-driven detection strategies share the same limitations, thus relaxing the stringent assumption of the knowledge of the system dynamics that is typically made in existing security works.

Future work includes the limitations of data-driven detection in stochastic, time-varying and nonlinear systems.

REFERENCES

- [1] A. A. Cárdenas, S. Amin, and S. S. Sastry. Research challenges for the security of control systems. In *Proceedings of the 3rd Conference on Hot Topics in Security*, pages 6:1–6:6, Berkeley, CA, USA, 2008.
- [2] F. Pasqualetti, F. Dörfler, and F. Bullo. Attack detection and identification in cyber-physical systems. *IEEE Transactions on Automatic Control*, 58(11):2715–2729, 2013.
- [3] Y. Chen, S. Kar, and J. M. F. Moura. Dynamic attack detection in cyber-physical systems with side initial state information. *IEEE Transactions on Automatic Control*, 62(9):4618–4624, 2016.
- [4] S. Papadimitriou, J. Sun, and C. Faloutsos. Streaming pattern discovery in multiple time-series. In *Proc. of the 31st International Conference on Very Large Data Bases*, pages 697–708, 2005.
- [5] D. Shi, Z. Guo, K. H. Johansson, and L. Shi. Causality countermeasures for anomaly detection in cyber-physical systems. *IEEE Transactions on Automatic Control*, 63(2):386–401, 2017.
- [6] R. Hink, J. M. Beaver, M. A. Buckner, T. Morris, U. Adhikari, and S. Pan. Machine learning for power system disturbance and cyber-attack discrimination. In *7th International Symposium on Resilient Control Systems*, pages 1–8, 2014.
- [7] M. Kravchik and A. Shabtai. Detecting cyber attacks in industrial control systems using convolutional neural networks. In *Proc. of the ACM Workshop on Cyber-Physical Systems Security and Privacy*, pages 72–83, 2018.

- [8] J. Inoue, Y. Yamagata, Y. Chen, C. M. Poskitt, and J. Sun. Anomaly detection for a water treatment system using unsupervised machine learning. In *IEEE International Conference on Data Mining Workshops*, pages 1058–1065, 2017.
- [9] J. Goh, S. Adepu, M. Tan, and Z. S. Lee. Anomaly detection in cyber physical systems using recurrent neural networks. In *IEEE 18th International Symposium on High Assurance Systems Engineering (HASE)*, pages 140–145, 2017.
- [10] J. C. Willems, P. Rapisarda, I. Markovsky, and B. L. M. De Moor. A note on persistency of excitation. *Systems & Control Letters*, 54(4):325–329, 2005.
- [11] H. J. Van Waarde, J. Eising, H. L. Trentelman, and M. K. Camlibel. Data informativity: a new perspective on data-driven analysis and control. *IEEE Transactions on Automatic Control*, 2020.
- [12] C. De Persis and P. Tesi. Formulas for data-driven control: Stabilization, optimality and robustness. *IEEE Transactions on Automatic Control*, 65(3):909–924, 2020.
- [13] G. Baggio, V. Katewa, and F. Pasqualetti. Data-driven minimum-energy controls for linear systems. *IEEE Control Systems Letters*, 3(3):589–594, 2019.
- [14] Y. Liu, P. Ning, and M. K. Reiter. False data injection attacks against state estimation in electric power grids. *ACM Transactions on Information and System Security*, 14(1):1–33, 2011.
- [15] Y. Mo and B. Sinopoli. False data injection attacks in electricity markets. In *IEEE Int. Conf. on Smart Grid Communications*, pages 226–231, Gaithersburg, MD, October 2010.
- [16] R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, 1985.
- [17] M. Sain and J. Massey. Invertibility of linear time-invariant dynamical systems. *IEEE Transactions on Automatic Control*, 14(2):141–149, 1969.
- [18] J. H. Tu, C. W. Rowley, D. M. Luchtenburg, S. L. Brunton, and J. N. Kutz. On dynamic mode decomposition: Theory and applications. *Journal of Computational Dynamics*, 1(2):391–421, 2014.
- [19] J-N. Juang and R. S. Pappa. An eigensystem realization algorithm for modal parameter identification and model reduction. *Journal of Guidance, Control and Dynamics*, 8(5):620–627, 1985.

APPENDIX I

PROOF OF THEOREM 2.1

We have $\mathbf{w}(k+1) = S^{(m,\text{nom})\top} \mathcal{O}_N x(k+1) = S^{(m,\text{nom})\top} \mathcal{O}_N A x(k)$. Suppose $x(k) \notin \text{Ker}(\mathcal{O}_N)$, we get that $0 \neq \mathcal{O}_N x(k) \in \text{Col}(\mathcal{O}_N) = \text{Col}(S^{(m,\text{nom})})$, and we can express $\mathcal{O}_N x(k)$ as a linear combination of the columns of $S^{(m,\text{nom})}$. It then follows that $S^{(m,\text{nom})\top} \mathcal{O}_N x(k) \neq 0$. Moreover, if $x(k) \in \text{Ker}(\mathcal{O}_N)$, then we have $S^{(m,\text{nom})\top} \mathcal{O}_N x(k) = 0$. Therefore, we get that $\text{Ker}(S^{(m,\text{nom})\top} \mathcal{O}_N) = \text{Ker}(\mathcal{O}_N)$.

We note that $\mathbf{w}(k) \in \text{Col}(S^{(m,\text{nom})\top} \mathcal{O}_N) \cong \text{Ker}^\perp(S^{(m,\text{nom})\top} \mathcal{O}_N)$ for all $k \in \mathbb{N}$. We now express $x(k) \in \mathbb{R}^n = \text{Ker}(S^{(m,\text{nom})\top} \mathcal{O}_N) \oplus \text{Ker}^\perp(S^{(m,\text{nom})\top} \mathcal{O}_N)$ as $x(k) = \alpha(k) + \beta(k)$, where $\alpha(k) \in \text{Ker}(S^{(m,\text{nom})\top} \mathcal{O}_N)$ and $\beta(k) \in \text{Ker}^\perp(S^{(m,\text{nom})\top} \mathcal{O}_N)$. We have $\mathbf{w}(k) = S^{(m,\text{nom})\top} \mathcal{O}_N x(k) = S^{(m,\text{nom})\top} \mathcal{O}_N \beta(k)$. Since $S^{(m,\text{nom})\top} \mathcal{O}_N$ is a bijection from $\text{Ker}^\perp(S^{(m,\text{nom})\top} \mathcal{O}_N)$ to $\text{Col}(S^{(m,\text{nom})\top} \mathcal{O}_N)$, with $(S^{(m,\text{nom})\top} \mathcal{O}_N)^\dagger$ as its inverse, we get that $\beta(k) = (S^{(m,\text{nom})\top} \mathcal{O}_N)^\dagger \mathbf{w}(k)$.

Now, we have $\mathbf{w}(k+1) = S^{(m,\text{nom})\top} \mathcal{O}_N x(k+1) = S^{(m,\text{nom})\top} \mathcal{O}_N A x(k) =$

$$\begin{aligned} S^{(m,\text{nom})\top} \mathcal{O}_N A (\alpha(k) + \beta(k)) &= S^{(m,\text{nom})\top} \mathcal{O}_N A \alpha(k) + \\ S^{(m,\text{nom})\top} \mathcal{O}_N A \beta(k) &= S^{(m,\text{nom})\top} \mathcal{O}_N A \alpha(k) + \\ S^{(m,\text{nom})\top} \mathcal{O}_N A (S^{(m,\text{nom})\top} \mathcal{O}_N)^\dagger \mathbf{w}(k). & \quad \text{Since} \end{aligned}$$

$\alpha(k) \in \text{Ker}(S^{(m,\text{nom})\top} \mathcal{O}_N) = \text{Ker}(\mathcal{O}_N)$, we have $CA^i \alpha(k) = 0$ for all $i \in \{0, \dots, N-1\}$. Now, if $N \geq \nu$, the observability index of the system (1), then we have $\text{Rank}(\mathcal{O}_N) = \text{Rank}(\mathcal{O}_{N+1})$, and since $\text{Ker}(\mathcal{O}_{N+1}) \subseteq \text{Ker}(\mathcal{O}_N)$, it follows from the Rank-Nullity Theorem that $\text{Ker}(\mathcal{O}_N) = \text{Ker}(\mathcal{O}_{N+1})$. We therefore get that $CA^N \alpha(k) = 0$, which implies that $\mathcal{O}_N A \alpha(k) = 0$, or in other words, $A \alpha(k) \in \text{Ker}(\mathcal{O}_N) = \text{Ker}(S^{(m,\text{nom})\top} \mathcal{O}_N)$.

This is equivalent to stating that the $\text{Ker}(S^{(m,\text{nom})\top} \mathcal{O}_N)$ is an invariant subspace of A when $N \geq \nu$, the observability index of (1). Therefore, we get that $\mathbf{w}(k+1) = (S^{(m,\text{nom})\top} \mathcal{O}_N)^\dagger A (S^{(m,\text{nom})\top} \mathcal{O}_N)^\dagger \mathbf{w}(k)$.

APPENDIX II

PROOF OF THEOREM 2.4

We first recall that $Y_{N,T} = \mathcal{O}_N [x(0) \quad Ax(0) \quad \dots \quad A^{T-N} x(0)] = [\mathcal{O}_N x(0) \quad \mathcal{O}_N A x(0) \quad \dots \quad \mathcal{O}_N A^{T-N} x(0)]$. We have $\mathbf{w}(k) = S^\top \mathcal{O}_N x(k)$ for all $k \in \mathbb{N}$. Let $x(k) = \alpha(k) + \beta(k)$, where $\alpha(k) \in \text{Ker}(S^\top \mathcal{O}_N)$ and $\beta(k) \in \text{Ker}^\perp(S^\top \mathcal{O}_N)$. We therefore get that $\mathbf{w}(k) = S^\top \mathcal{O}_N \beta(k)$. Since $S^\top \mathcal{O}_N$ is a bijection from $\text{Ker}^\perp(S^\top \mathcal{O}_N)$ to $\text{Col}(S^\top \mathcal{O}_N)$, we have that $\beta(k) = (S^\top \mathcal{O}_N)^\dagger \mathbf{w}(k)$.

Now, since $N \geq \nu$ and $T - N \geq \mu - 1$, we get from Theorem 2.3 that $Y_{N,T}$ is of maximum rank and we have $\mathbf{w}(k+1) = S^\top \mathcal{O}_N x(k+1) = S^\top \mathcal{O}_N A x(k) = S^\top \mathcal{O}_N A (\alpha(k) + \beta(k)) = S^\top \mathcal{O}_N A \alpha(k) + S^\top \mathcal{O}_N A (S^\top \mathcal{O}_N)^\dagger \mathbf{w}(k)$.

Suppose $\alpha(k) \in \text{Ker}(\mathcal{O}_N)$, and since $N \geq \nu$ we get that $\text{Ker}(\mathcal{O}_N)$ is A -invariant, and we get that $A \alpha(k) \in \text{Ker}(\mathcal{O}_N)$, or in other words, $S^\top \mathcal{O}_N A \alpha(k) = 0$. Now suppose, on the other hand, that $\alpha(k) \in \text{Ker}^\perp(\mathcal{O}_N) \cap \text{Ker}(S^\top \mathcal{O}_N)$, we would then indeed have $x(k) \in \text{Ker}^\perp(\mathcal{O}_N)$. Moreover, let $V \subseteq \text{Ker}^\perp(\mathcal{O}_N)$ be the smallest A -invariant subset of $\text{Ker}^\perp(\mathcal{O}_N)$ containing $x(0)$. This implies that $x(k) \in V \subseteq \text{Ker}^\perp(\mathcal{O}_N)$. Since the columns of S form a basis of $\mathcal{O}_N V$, we get that $\mathcal{O}_N x(k)$ can be expressed as a linear combination of the columns of S , i.e., that $\mathcal{O}_N x(k) = S v(k)$, with $v(k) \neq 0$ and we get $x(k) = \mathcal{O}_N^\dagger S v(k)$. It then follows that $\mathcal{O}_N x(k)$ does not have a component orthogonal to all the column vectors of S , which implies that $\alpha(k) = 0$ and we get a contradiction. Therefore, $\alpha(k) \notin \text{Ker}^\perp(\mathcal{O}_N) \cap \text{Ker}(S^\top \mathcal{O}_N)$. Thus, we get that $\mathbf{w}(k+1) = (S^\top \mathcal{O}_N)^\dagger A (S^\top \mathcal{O}_N)^\dagger \mathbf{w}(k)$, and we can write $\vec{W} = (S^\top \mathcal{O}_N)^\dagger A (S^\top \mathcal{O}_N)^\dagger W$. It is clear from the above that $\text{Col}(\vec{W}^\top) \subseteq \text{Col}(W^\top)$, from which we infer that $(S^\top \mathcal{O}_N)^\dagger A (S^\top \mathcal{O}_N)^\dagger = \vec{W} W^\dagger$ is the global minimizer for the minimization problem (6) and that the minimum value is zero.