

# Defending Surveillance Sensor Networks Against Data-Injection Attacks via Trusted Nodes

Roberto López-Valcarce  
University of Vigo, Spain  
valcarce@gts.uvigo.es

Daniel Romero  
University of Agder, Norway  
daniel.romero@uia.no

**Abstract**—By injecting false data through compromised sensors, an adversary can drive the probability of detection in a sensor network-based spatial field surveillance system to arbitrarily low values. As a countermeasure, a small subset of sensors may be secured. Leveraging the theory of Matched Subspace Detection, we propose and evaluate several detectors that add robustness to attacks when such trusted nodes are available. Our results reveal the performance-security tradeoff of these schemes and can be used to determine the number of trusted nodes required for a given performance target.

**Index Terms**—Adversarial signal processing, Byzantine sensors, cyber security, sensor networks.

## I. INTRODUCTION

Sensor networks for surveillance and environmental monitoring [1], [2] comprise a large number of nodes measuring some physical phenomena and reporting to a fusion center (FC). Unattended nodes in typical deployments over large areas are susceptible to external attacks [3], [4], e.g. an adversary may capture some nodes and change the content of data packets, or directly alter the environment around some sensors. Standard cryptographic measures are ineffective against such data-injection attacks.

We focus on typical dense deployments, for which measurements of the monitored physical phenomenon exhibit spatial smoothness [5]. This allows for parsimonious parametric modeling of the corresponding spatial field, which can be exploited for inference purposes [6]–[8]. As in [9], a linear model in spatially independent Gaussian noise with unknown variance is adopted. The FC tests for the presence of a spatially smooth field (bacterial activity, toxic chemical spill, radioactivity, etc.). A related model with data-injection attacks for state *estimation* (rather than *detection*) in cyber-physical systems has been used in [10], [11]. An adversary may inject false data in a number of sensors to preclude the FC from detecting the spatial field. If the number of compromised sensors is strictly larger than the signal subspace dimension, the adversary will be able to maximize system degradation [9]. Since the effects can be disastrous, it is important to devise adequate defense mechanisms against such attacks.

Work funded by the Agencia Estatal de Investigación (Spain) and the European Regional Development Fund (ERDF) (projects TEC2013-47020-C2-1-R, TEC2015-69648-REDC, TEC2016-76409-C2-2-R), and by the Xunta de Galicia and ERDF (projects GRC2013/009, R2014/037 and ED431G/04 Agrupación Estratégica Consolidada de Galicia accreditation 2016-2019).

To this end, we propose to endow the network with *trusted nodes*. By investing additional resources, network deployment may be planned so that a number of sensors are placed at secure locations, out of reach to adversaries. In this extended form of *tamper resistance* [12], not only physical protection against unauthorized attempts to read or modify the content of the device must be provided, but also against intentional alteration of the measured field near the node. For example, in clustered network topologies [13], [14], cluster heads are natural candidates for becoming trusted nodes. Alternatively, the FC may obtain estimates of node reliability via reputation metrics based on previous interactions [15], [16], so that only those nodes with a sufficiently high reputation metric are labeled as trusted. Our goal is to exploit the availability of trusted nodes to make spatial field detection robust to malicious actions. To this end, a number of schemes are developed based on the framework of Matched Subspace Detection [17]. These schemes will be shown to exhibit different tradeoffs between robustness and performance.

*Notation:* For a matrix  $\mathbf{A}$ ,  $\mathbf{A}^\dagger$  denotes its pseudoinverse, and  $\mathcal{R}(\mathbf{A})$ ,  $\mathcal{R}^\perp(\mathbf{A})$  its column and left null spaces, respectively. The perpendicular projection matrices onto  $\mathcal{R}(\mathbf{A})$  and  $\mathcal{R}^\perp(\mathbf{A})$  are respectively denoted by  $\mathbf{P}_\mathbf{A} = \mathbf{A}\mathbf{A}^\dagger$  and  $\mathbf{P}_\mathbf{A}^\perp = \mathbf{I} - \mathbf{P}_\mathbf{A}$ . The Gaussian distribution with mean  $\boldsymbol{\mu}$  and covariance  $\mathbf{C}$  is denoted by  $\mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$ . By  $\mathcal{F}_{\nu_1, \nu_2}$ ,  $\mathcal{F}'_{\nu_1, \nu_2}(\lambda_1)$  and  $\mathcal{F}''_{\nu_1, \nu_2}(\lambda_1, \lambda_2)$  we respectively denote the central, non-central, and doubly noncentral  $F$ -distributions with  $\nu_1$  and  $\nu_2$  degrees of freedom (d.o.f.) and noncentrality parameters  $\lambda_1$ ,  $\lambda_2$  [18].

## II. PROBLEM STATEMENT AND PRELIMINARIES

Consider a surveillance network deployed to monitor the presence or absence of some spatial field, consisting of an FC collecting data from  $n$  sensors. These data are comprised in vector  $\mathbf{y} \triangleq [y_1 \cdots y_n]^T$ , modeled as

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{a} + \mathbf{w} \in \mathbb{R}^n, \quad (1)$$

where  $\mathbf{H} \in \mathbb{R}^{n \times m}$  is a known full-rank matrix,  $\mathbf{x} \in \mathbb{R}^m$  is related to the monitored physical process, and  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$  is measurement noise with unknown variance  $\sigma^2$ . The *attack vector*  $\mathbf{a} = [a_1 \cdots a_n]^T \in \mathbb{R}^n$  is injected by the adversary. The goal of the network is to detect the presence of the spatial field, i.e., to decide whether or not  $\mathbf{x} = \mathbf{0}$ . The basis expansion model  $\mathbf{H}\mathbf{x}$  of the signal component is fairly general; it hinges

on the assumption that the spatial field is sufficiently smooth so it can be parameterized by a low-dimensional  $\mathbf{x}$  ( $m \ll n$ ). It fits a wide range of signal representations based on Fourier or Discrete Cosine Transforms, polynomial bases, splines, etc. [6], [8]. Thus,  $\mathbf{H}$  depends on the chosen basis expansion and sensor locations.

After instigating some event (e.g., chemical spill) which will eventually generate a spatial field measurable by the network ( $\mathbf{x} \neq \mathbf{0}$ ), the adversary's goal is to prevent this event from being detected, by injecting false data into a subset of  $k \ll n$  compromised sensors. Thus,  $a_i$  is freely selected by the adversary if the  $i$ -th node was compromised; else,  $a_i = 0$ .

Our goal is to investigate defense mechanisms to make the network robust to these data-injection attacks. To this end, first we briefly review Matched Subspace Detection [17], since the problems studied in the sequel fit in this class. Subsequently, we will review the effect of false data injection in an attack-unaware network [9], which will motivate the search for alternatives robust to this kind of attacks.

### A. Matched Subspace Detectors

Consider the general model

$$\begin{aligned} \mathbf{y} &= \mathbf{A}\boldsymbol{\theta} + \mathbf{v}, & \mathbf{A} \in \mathbb{R}^{n \times p}, & \boldsymbol{\theta} \in \mathbb{R}^p, \\ \mathbf{v} &\sim \mathcal{N}(\mathbf{B}\boldsymbol{\phi}, \sigma^2 \mathbf{I}_n), & \mathbf{B} \in \mathbb{R}^{n \times q}, & \boldsymbol{\phi} \in \mathbb{R}^q, \end{aligned} \quad (2)$$

with  $\mathbf{A}$ ,  $\mathbf{B}$  known, and  $\boldsymbol{\theta}$ ,  $\boldsymbol{\phi}$ ,  $\sigma^2$  unknown. The matrices  $\mathbf{A}$ ,  $\mathbf{B}$  and  $[\mathbf{A} \ \mathbf{B}]$  are assumed full-rank. We wish to test for the presence of the signal component  $\mathbf{A}\boldsymbol{\theta}$ :  $\mathcal{H}_0 : \mathbf{y} \sim \mathcal{N}(\mathbf{B}\boldsymbol{\phi}, \sigma^2 \mathbf{I}_n)$  vs.  $\mathcal{H}_1 : \mathbf{y} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\theta} + \mathbf{B}\boldsymbol{\phi}, \sigma^2 \mathbf{I}_n)$ . The component  $\mathbf{B}\boldsymbol{\phi}$  can be regarded as an interference term. Thus, in this model, the signal term is known to lie in the subspace  $\mathcal{R}(\mathbf{A})$ , whereas the interference is known to lie in the subspace  $\mathcal{R}(\mathbf{B})$ .

Let  $\mathbf{G} = \mathbf{P}_B^\perp \mathbf{A}$ . Provided that  $p + q < n$ , the Generalized Likelihood Ratio Test (GLRT) exists and is given by

$$T \triangleq \frac{n-p-q}{p} \frac{\|\mathbf{P}_G \mathbf{P}_B^\perp \mathbf{y}\|^2}{\|\mathbf{P}_G^\perp \mathbf{P}_B^\perp \mathbf{y}\|^2} \underset{\mathcal{H}_0}{\overset{\mathcal{H}_1}{\gtrless}} \gamma, \quad (3)$$

and it is uniformly most powerful invariant (UMPI) [17]. In (3), the interference is first cancelled by  $\mathbf{P}_B^\perp$ , after which the signal component is known to lie in  $\mathcal{R}(\mathbf{G})$ . Then (3) compares the energies of  $\mathbf{P}_B^\perp \mathbf{y}$  in this subspace and in its orthogonal complement  $\mathcal{R}^\perp(\mathbf{G})$ . The test statistic  $T$  is  $F$ -distributed as  $T | \mathcal{H}_0 \sim \mathcal{F}_{p, n-p-q}$  and  $T | \mathcal{H}_1 \sim \mathcal{F}'_{p, n-p-q} \left( \frac{\|\mathbf{G}\boldsymbol{\theta}\|^2}{\sigma^2} \right)$ . For fixed  $\frac{\|\mathbf{G}\boldsymbol{\theta}\|^2}{\sigma^2}$ , detection performance improves as the first number of d.o.f. (signal subspace dimension  $p$ ) decreases, and as the second number of d.o.f.  $n-p-q$  increases.

### B. Field detection with an unprotected network

Consider now the original model (1), and suppose that the network is oblivious to potential attacks, meaning that it assumes  $\mathbf{a} = \mathbf{0}$ . Under this assumption, the problem of testing for field presence, i.e.,  $\mathcal{H}_0 : \mathbf{x} = \mathbf{0}$  vs.  $\mathcal{H}_1 : \mathbf{x} \neq \mathbf{0}$ , can be cast in the framework of (2) with  $\mathbf{A}\boldsymbol{\theta} = \mathbf{H}\mathbf{x}$  and  $\mathbf{B}\boldsymbol{\phi} = \mathbf{0}$  (and thus  $p = m$  and  $q = 0$ ). The corresponding GLRT is

$$T_F \triangleq \frac{n-m}{m} \frac{\|\mathbf{P}_H \mathbf{y}\|^2}{\|\mathbf{P}_H^\perp \mathbf{y}\|^2} \underset{\mathcal{H}_0}{\overset{\mathcal{H}_1}{\gtrless}} \gamma_F, \quad (4)$$

and in the absence of attacks, one has  $T_F | \mathcal{H}_0 \sim \mathcal{F}_{m, n-m}$  and  $T_F | \mathcal{H}_1 \sim \mathcal{F}'_{m, n-m}(\rho)$ , where  $\rho \triangleq \frac{\|\mathbf{H}\mathbf{x}\|^2}{\sigma^2}$  is the Signal-to-Noise Ratio (SNR). However, as shown in [9], the presence of an attack changes the distribution of  $T_F$  under  $\mathcal{H}_1$  to

$$T_F \sim \mathcal{F}''_{m, n-m}(\rho_{||}, \rho_{\perp}) \quad \text{with} \quad \begin{cases} \rho_{||} \triangleq \frac{\|\mathbf{H}\mathbf{x} + \mathbf{P}_H \mathbf{a}\|^2}{\sigma^2}, \\ \rho_{\perp} \triangleq \frac{\|\mathbf{P}_H^\perp \mathbf{a}\|^2}{\sigma^2}. \end{cases} \quad (5)$$

Under (5), the detection probability of (4) is monotonically increasing in  $\rho_{||}$  and decreasing in  $\rho_{\perp}$ . As in [9], we assume that the adversary knows  $\mathbf{H}$  (which could be inferred from sensor locations), but not  $\mathbf{x}$  (which depends on unknown environmental variables affecting field diffusion). Thus, to avoid accidentally increasing  $\rho_{||}$ , the adversary may choose an *orthogonal attack*  $\mathbf{a} \in \mathcal{R}^\perp(\mathbf{H})$ ; for this, it suffices to capture  $k > m$  sensors [9]. The performance of (4) can be drastically degraded by this kind of attack, in the sense that its detection probability can be made arbitrarily small by increasing the attack power  $\|\mathbf{a}\|^2$  [9]. Clearly, appropriate countermeasures are needed.

Note that directly applying the results in Sec. II-A to (1) (by regarding  $\mathbf{a}$  as the interference term) is not possible: the subset of compromised sensors is unknown to the FC, so it does not know any subspace in which  $\mathbf{a}$  must lie. At most, by the reasoning in Sec. II-B, one could assume that  $\mathbf{a} \in \mathcal{R}^\perp(\mathbf{H})$ , but then, in the framework of Sec. II-A,  $p = m$  and  $q = n - m$ : this violates the constraint  $p + q < n$ , so the GLRT does not exist. As shown in the sequel, this situation changes once trusted nodes become available.

## III. FIELD DETECTION WITH TRUSTED NODES

One sensible way to provide robustness to attacks is to secure a subset of nodes. Thus, suppose that a number  $t$  of them, say nodes  $1, \dots, t$ , are trusted. Then the attack vector must be of the form  $\mathbf{a} = [\mathbf{0}_t^T \ \tilde{\mathbf{a}}^T]^T$  with  $\tilde{\mathbf{a}} \in \mathbb{R}^{n-t}$ . Letting  $\tilde{\mathbf{H}} \in \mathbb{R}^{t \times m}$  and  $\tilde{\mathbf{H}} \in \mathbb{R}^{(n-t) \times m}$  comprise the first  $t$  and last  $n-t$  rows of  $\mathbf{H}$  respectively, model (1) can be split as

$$\bar{\mathbf{y}} = \tilde{\mathbf{H}}\mathbf{x} + \tilde{\mathbf{w}}, \quad \tilde{\mathbf{y}} = \tilde{\mathbf{H}}\mathbf{x} + \tilde{\mathbf{a}} + \tilde{\mathbf{w}}, \quad (6)$$

with  $\bar{\mathbf{y}} \triangleq [y_1 \cdots y_t]^T$ ,  $\tilde{\mathbf{y}} \triangleq [y_{t+1} \cdots y_n]^T$  (data from trusted and untrusted nodes respectively), and  $\tilde{\mathbf{w}}$ ,  $\tilde{\mathbf{w}}$  the corresponding noise vectors. The availability of trusted nodes constrains the attack vector to a certain subspace, allowing the application of the theory from Sec. II-A. Now, two strategies are possible: the first one makes no assumptions on the attack vector  $\tilde{\mathbf{a}}$ , whereas the second assumes a structured (orthogonal) attack.

### A. Unstructured attack

With no assumptions on  $\tilde{\mathbf{a}}$ , then in the framework of Sec. II-A, one may take

$$\mathbf{A} = \mathbf{H}, \quad \mathbf{B} = \begin{bmatrix} \mathbf{0} \\ \mathbf{I}_{n-t} \end{bmatrix} \Rightarrow \begin{cases} p = m, \\ q = n-t. \end{cases} \quad (7)$$

Hence, to have  $p + q < n$  one needs  $t > m$ , i.e., the number of trusted nodes must exceed the dimension of the signal

subspace. Then, the GLRT statistic for field detection is given by

$$T_{\text{FU}} \triangleq \frac{t-m}{m} \frac{\|\mathbf{P}_{\tilde{\mathbf{H}}}\tilde{\mathbf{y}}\|^2}{\|\mathbf{P}_{\tilde{\mathbf{H}}}^\perp\tilde{\mathbf{y}}\|^2} \sim \mathcal{F}'_{m,t-m} \left( \frac{\|\tilde{\mathbf{H}}\mathbf{x}\|^2}{\sigma^2} \right). \quad (8)$$

Attacks have no effect on the performance of (8), but all data from untrusted nodes is discarded, which is clearly inefficient. Alternatives exhibiting a better robustness-performance trade-off will be investigated next.

### B. Structured attack

If an orthogonal attack  $\mathbf{a} \in \mathcal{R}^\perp(\mathbf{H})$  is assumed, then the resulting interference subspace can be further constrained. This subspace can be characterized as follows.

**Lemma 1.** *If  $\mathbf{a} = [\mathbf{0}_t^T \tilde{\mathbf{a}}^T]^T$  and  $\mathbf{a} \in \mathcal{R}^\perp(\mathbf{H})$ , then  $\tilde{\mathbf{a}} \in \mathcal{R}^\perp(\tilde{\mathbf{H}})$ .*

**Proof.**  $\mathbf{a} \in \mathcal{R}^\perp(\mathbf{H})$  means  $\mathbf{H}^T \mathbf{a} = \mathbf{0}$ . Since  $\mathbf{a} = [\mathbf{0}_t^T \tilde{\mathbf{a}}^T]^T$ , one has  $\mathbf{H}^T \mathbf{a} = \mathbf{0} \Rightarrow \tilde{\mathbf{H}}^T \tilde{\mathbf{a}} = \mathbf{0}$ , hence  $\tilde{\mathbf{a}} \in \mathcal{R}^\perp(\tilde{\mathbf{H}})$ . ■

Let  $\tilde{r} = \text{rank } \tilde{\mathbf{H}}$ , and let the columns of  $\tilde{\mathbf{U}}_\perp \in \mathbb{R}^{(n-t) \times (n-t-\tilde{r})}$  be an orthonormal basis for  $\mathcal{R}^\perp(\tilde{\mathbf{H}})$ , so that  $\mathbf{P}_{\tilde{\mathbf{H}}}^\perp = \tilde{\mathbf{U}}_\perp \tilde{\mathbf{U}}_\perp^T$ . Assume that  $\mathbf{a} \in \mathcal{R}^\perp(\mathbf{H})$ . Then, in view of Lemma 1, field detection can be cast in the framework of Sec. II-A by setting

$$\mathbf{A} = \mathbf{H}, \quad \mathbf{B} = \begin{bmatrix} \mathbf{0} \\ \tilde{\mathbf{U}}_\perp \end{bmatrix} \Rightarrow \begin{cases} p = m, \\ q = n-t-\tilde{r}. \end{cases} \quad (9)$$

The GLRT exists if  $p+q < n$ , i.e.,  $t > m - \tilde{r}$  (in particular, if the matrix  $\tilde{\mathbf{H}}$  has full column rank  $\tilde{r} = m$ , then the GLRT exists for any number of trusted nodes  $t > 0$ ), in which case

$$\left. \begin{aligned} \mathbf{P}_B^\perp &= \mathbf{I} - \mathbf{P} \\ \mathbf{G} &= (\mathbf{I} - \mathbf{P})\mathbf{H} = \mathbf{H} \end{aligned} \right\} \text{ with } \mathbf{P} = \begin{bmatrix} \mathbf{0} & \\ & \mathbf{P}_{\tilde{\mathbf{H}}}^\perp \end{bmatrix}. \quad (10)$$

Since  $\mathbf{P}_H \mathbf{P} = \mathbf{0}$ , it follows that  $\mathbf{P}_G \mathbf{P}_B^\perp = \mathbf{P}_H$  and  $\mathbf{P}_G^\perp \mathbf{P}_B^\perp = \mathbf{P}_{\tilde{\mathbf{H}}}^\perp - \mathbf{P}$ . This results in the GLRT statistic

$$T_{\text{FS}} \triangleq \frac{t+\tilde{r}-m}{m} \frac{\|\mathbf{P}_H \mathbf{y}\|^2}{\|\mathbf{P}_{\tilde{\mathbf{H}}}^\perp \mathbf{y}\|^2 - \|\mathbf{P}_{\tilde{\mathbf{H}}}^\perp \tilde{\mathbf{y}}\|^2} \quad (11)$$

$$\sim \mathcal{F}'_{m,t+\tilde{r}-m} \left( \frac{\|\mathbf{H}\mathbf{x}\|^2}{\sigma^2} \right), \quad (12)$$

which coincides with the distribution of (4) under no attacks,  $T_{\text{F}} \sim \mathcal{F}'_{m,n-m} \left( \frac{\|\mathbf{H}\mathbf{x}\|^2}{\sigma^2} \right)$ , if  $n$  is replaced by  $t + \tilde{r}$ . The operational SNR of both tests is the same, in contrast with (8).

To analyze the effect of an attack not complying with the original assumption  $\mathbf{a} \in \mathcal{R}^\perp(\mathbf{H})$ , note from Lemma 1 that this implies  $\mathbf{a} = [\mathbf{0}_t^T \tilde{\mathbf{a}}^T]^T$  with  $\tilde{\mathbf{a}} \notin \mathcal{R}^\perp(\tilde{\mathbf{H}})$ . Hence, let

$$\mathbf{a} = \underbrace{\begin{bmatrix} \mathbf{0}_t \\ \tilde{\mathbf{a}}_{\parallel} \end{bmatrix}}_{=\mathbf{a}_{\parallel}} + \underbrace{\begin{bmatrix} \mathbf{0}_t \\ \tilde{\mathbf{a}}_{\perp} \end{bmatrix}}_{=\mathbf{a}_{\perp}} \text{ with } \begin{cases} \tilde{\mathbf{a}}_{\parallel} \in \mathcal{R}(\tilde{\mathbf{H}}), \\ \tilde{\mathbf{a}}_{\perp} \in \mathcal{R}^\perp(\tilde{\mathbf{H}}). \end{cases} \quad (13)$$

(Note that, in general,  $\mathbf{a}_{\parallel}$  does not lie either in  $\mathcal{R}(\mathbf{H})$  nor in  $\mathcal{R}^\perp(\mathbf{H})$ ). Then the distribution of  $T_{\text{FS}}$  changes from (12) to

$$T_{\text{FS}} \sim \mathcal{F}''_{m,t+\tilde{r}-m} \left( \frac{\|\mathbf{H}\mathbf{x} + \mathbf{P}_H \mathbf{a}_{\parallel}\|^2}{\sigma^2}, \frac{\|\mathbf{P}_{\tilde{\mathbf{H}}}^\perp \mathbf{a}_{\parallel}\|^2}{\sigma^2} \right). \quad (14)$$

The detection probability is monotonically increasing (resp. decreasing) in the first (resp. second) noncentrality parameter of (14). Again, since  $\mathbf{x}$  is unknown to the adversary, he is unable to design the attack to reduce  $\|\mathbf{H}\mathbf{x} + \mathbf{P}_H \mathbf{a}_{\parallel}\|^2$ ; with a "poor" choice for  $\mathbf{a}_{\parallel}$ , this term may actually increase. Thus, a sensible approach for the attacker would be to maximize  $\|\mathbf{P}_{\tilde{\mathbf{H}}}^\perp \mathbf{a}_{\parallel}\|^2$  subject to  $\mathbf{P}_H \mathbf{a}_{\parallel} = \mathbf{0}$ . However, one has:

**Lemma 2.** *If  $\mathbf{a}_{\parallel} = [\mathbf{0}_t^T \tilde{\mathbf{a}}_{\parallel}^T]^T$  with  $\tilde{\mathbf{a}}_{\parallel} \in \mathcal{R}(\tilde{\mathbf{H}})$ , then  $\mathbf{P}_H \mathbf{a}_{\parallel} = \mathbf{0}$  implies  $\mathbf{a}_{\parallel} = \mathbf{0}$ .*

**Proof.**  $\mathbf{P}_H \mathbf{a}_{\parallel} = \mathbf{0}$  means  $\mathbf{a}_{\parallel} \in \mathcal{R}^\perp(\mathbf{H})$ , i.e.,  $\mathbf{H}^T \mathbf{a}_{\parallel} = \tilde{\mathbf{H}}^T \tilde{\mathbf{a}}_{\parallel} = \mathbf{0}$ . Thus,  $\tilde{\mathbf{a}}_{\parallel} \in \mathcal{R}^\perp(\tilde{\mathbf{H}})$ , but since  $\tilde{\mathbf{a}}_{\parallel} \in \mathcal{R}(\tilde{\mathbf{H}})$  by assumption, it follows that  $\tilde{\mathbf{a}}_{\parallel} = \mathbf{0}$ , hence  $\mathbf{a}_{\parallel} = \mathbf{0}$ . ■

Thus, any attack with  $\|\mathbf{P}_{\tilde{\mathbf{H}}}^\perp \mathbf{a}_{\parallel}\|^2 > 0$  must yield  $\mathbf{P}_H \mathbf{a}_{\parallel} \neq \mathbf{0}$ , so the effect on detection probability cannot be determined a priori: it may very well increase due to the attack. This establishes the robustness of (11) to data-injection attacks.

## IV. JOINT FIELD-ATTACK DETECTION

The schemes in Sec. III attempt to block  $\mathbf{a}$  before testing for field presence. Noting that the adversary is only interested in launching an attack if a field is present suggests an alternative approach, by which the FC jointly tests for the presence of the spatial field *or* the attack vector, i.e.,  $\mathcal{H}_0 : \mathbf{H}\mathbf{x} + \mathbf{a} = \mathbf{0}$  vs.  $\mathcal{H}_1 : \mathbf{H}\mathbf{x} + \mathbf{a} \neq \mathbf{0}$ . Again, nodes  $1, \dots, t$  are assumed trusted, so that  $\mathbf{a} = [\mathbf{0}_t^T \tilde{\mathbf{a}}^T]^T$ .

### A. Unstructured attack

With no assumptions on  $\tilde{\mathbf{a}}$ , then in the framework of Sec. II-A,  $\mathbf{B}$  is empty (thus,  $\mathbf{P}_B^\perp = \mathbf{I}$ ) and  $\mathbf{A}\boldsymbol{\theta} = \mathbf{H}\mathbf{x} + \mathbf{a}$ . Hence,

$$\mathbf{A} = \begin{bmatrix} \tilde{\mathbf{H}} & \mathbf{0} \\ \tilde{\mathbf{H}} & \mathbf{I}_{n-t} \end{bmatrix} \Rightarrow \begin{cases} p = m+n-t, \\ q = 0. \end{cases} \quad (15)$$

The GLRT exists if  $t > m$ , and its test statistic is given by

$$T_{\text{JU}} \triangleq \frac{t-m}{m+n-t} \frac{\|\mathbf{P}_{\tilde{\mathbf{H}}}\tilde{\mathbf{y}}\|^2 + \|\tilde{\mathbf{y}}\|^2}{\|\mathbf{P}_{\tilde{\mathbf{H}}}^\perp \tilde{\mathbf{y}}\|^2} \quad (16)$$

$$\sim \mathcal{F}'_{m+n-t,t-m} \left( \frac{\|\mathbf{H}\mathbf{x} + \mathbf{a}\|^2}{\sigma^2} \right). \quad (17)$$

### B. Structured attack

Assuming an orthogonal attack  $\mathbf{a} \in \mathcal{R}^\perp(\mathbf{H})$  (and therefore  $\tilde{\mathbf{a}} \in \mathcal{R}^\perp(\tilde{\mathbf{H}})$ , in view of Lemma 1), we again set  $\mathbf{B}$  to be empty, whereas now

$$\mathbf{A} = \begin{bmatrix} \tilde{\mathbf{H}} & \mathbf{0} \\ \tilde{\mathbf{H}} & \tilde{\mathbf{U}}_\perp \end{bmatrix} \Rightarrow \begin{cases} p = m+n-t-\tilde{r}, \\ q = 0. \end{cases} \quad (18)$$

Then  $\mathbf{G} = \mathbf{A}$  and, with  $\mathbf{P}$  as in (10), one has  $\mathbf{P}_G = \mathbf{P}_H + \mathbf{P}$  and  $\mathbf{P}_G^\perp = \mathbf{P}_{\tilde{\mathbf{H}}}^\perp - \mathbf{P}$ . The GLRT exists if  $t > m - \tilde{r}$  and, since  $\mathbf{P}_H \mathbf{P} = \mathbf{0}$ , its statistic can be written as

$$T_{\text{JS}} \triangleq \frac{t+\tilde{r}-m}{m+n-t-\tilde{r}} \frac{\|\mathbf{P}_H \mathbf{y}\|^2 + \|\mathbf{P}_{\tilde{\mathbf{H}}}^\perp \tilde{\mathbf{y}}\|^2}{\|\mathbf{P}_{\tilde{\mathbf{H}}}^\perp \mathbf{y}\|^2 - \|\mathbf{P}_{\tilde{\mathbf{H}}}^\perp \tilde{\mathbf{y}}\|^2} \quad (19)$$

$$\sim \mathcal{F}'_{m+n-t-\tilde{r},t+\tilde{r}-m} \left( \frac{\|\mathbf{H}\mathbf{x}\|^2}{\sigma^2} + \frac{\|\mathbf{a}\|^2}{\sigma^2} \right) \quad (20)$$

with (20) valid if  $\mathbf{a} \in \mathcal{R}^\perp(\mathbf{H})$ . Otherwise, with  $\mathbf{a}$  given by (13), the distribution of  $T_{\text{JS}}$  changes to

$$\mathcal{F}''_{m+n-t-\tilde{r}, t+\tilde{r}-m} \left( \frac{\|\mathbf{H}\mathbf{x} + \mathbf{P}_H \mathbf{a}_\parallel\|^2}{\sigma^2} + \frac{\|\mathbf{a}_\perp\|^2}{\sigma^2}, \frac{\|\mathbf{P}_H^\perp \mathbf{a}_\parallel\|^2}{\sigma^2} \right) \quad (21)$$

Any  $\mathbf{a}_\perp \neq \mathbf{0}$  will contribute to increasing the probability of detection. Regarding the component  $\mathbf{a}_\parallel$ , in view of Lemma 2, the adversary cannot increase the term  $\|\mathbf{P}_H^\perp \mathbf{a}_\parallel\|^2$  in (21) without the risk of increasing the term  $\|\mathbf{H}\mathbf{x} + \mathbf{P}_H \mathbf{a}_\parallel\|^2$ .

### C. OR detector

In this approach, the field detection test (4), which assumes attacks are absent, is applied simultaneously with a second detector testing for the presence of an attack, as described below, and then the results of both tests are fused using the OR rule. The rationale for this is that, whereas the performance of the field detector (4) will degrade as the attack power increases, that of the attack detector should improve. Since the adversary will only launch an attack to cover up the presence of the spatial field, an OR fusion rule is well motivated.

The detector for  $\mathcal{H}_0 : \mathbf{a} = \mathbf{0}$  vs.  $\mathcal{H}_1 : \mathbf{a} \neq \mathbf{0}$  can be developed as per Sec. II-A, by regarding the attack as signal and the field as interference. We assume an orthogonal attack (a similar approach is possible for unstructured attacks). Thus,

$$\mathbf{A} = \begin{bmatrix} \mathbf{0} \\ \tilde{\mathbf{U}}_\perp \end{bmatrix}, \quad \mathbf{B} = \mathbf{H} \Rightarrow \begin{cases} p = n - t - \tilde{r}, \\ q = m. \end{cases} \quad (22)$$

The GLRT exists if  $t > m - \tilde{r}$ . Since  $\mathbf{P}_B^\perp = \mathbf{P}_H^\perp$ , one has  $\mathbf{G} = \mathbf{P}_H^\perp \mathbf{A} = \mathbf{A}$ , which has orthonormal columns; hence,  $\mathbf{P}_G = \mathbf{A}\mathbf{A}^T = \mathbf{P}$ , with  $\mathbf{P}$  given by (10). Then  $\mathbf{P}_G \mathbf{P}_B^\perp = \mathbf{P} \mathbf{P}_H^\perp = \mathbf{P}$ , and  $\mathbf{P}_G^\perp \mathbf{P}_B^\perp = \mathbf{P}_H^\perp - \mathbf{P}$ , yielding the test statistic

$$T_{\text{AS}} \triangleq \frac{t + \tilde{r} - m}{n - t - \tilde{r}} \frac{\|\mathbf{P}_H^\perp \tilde{\mathbf{y}}\|^2}{\|\mathbf{P}_H^\perp \tilde{\mathbf{y}}\|^2 - \|\mathbf{P}_H^\perp \tilde{\mathbf{y}}\|^2} \quad (23)$$

$$\sim \mathcal{F}'_{n-t-\tilde{r}, t+\tilde{r}-m} \left( \frac{\|\mathbf{a}\|^2}{\sigma^2} \right), \quad (24)$$

with (24) valid if  $\mathbf{a} \in \mathcal{R}^\perp(\mathbf{H})$ . Otherwise, let  $\mathbf{a}$  be given by (13). Then, the distribution of  $T_{\text{AS}}$  changes to

$$T_{\text{AS}} \sim \mathcal{F}''_{n-t-\tilde{r}, t+\tilde{r}-m} \left( \frac{\|\mathbf{a}_\perp\|^2}{\sigma^2}, \frac{\|\mathbf{P}_H^\perp \mathbf{a}_\parallel\|^2}{\sigma^2} \right). \quad (25)$$

The detection performance of (25) degrades with  $\|\mathbf{P}_H^\perp \mathbf{a}_\parallel\|^2$ , but again by virtue of Lemma 2, the adversary cannot increase such term without the risk of increasing the term  $\rho_\parallel$  in (5) and, in turn, the detection probability of the spatial field detector.

The following result holds now:

**Lemma 3.** *In the absence of attacks, the test statistics  $T_{\text{F}}$  from (4) and  $T_{\text{AS}}$  from (23) are statistically independent.*

We skip the proof for lack of space. By Lemma 3, the probability of false alarm of the OR detector is  $P_{\text{FA}} = P_{\text{FA},1} + P_{\text{FA},2} - P_{\text{FA},1}P_{\text{FA},2}$ , with  $P_{\text{FA},i}$  the probabilities of false alarm of the individual tests. A degree of freedom is available to choose  $P_{\text{FA},i}$ ; for a target  $P_{\text{FA}}$ , we suggest

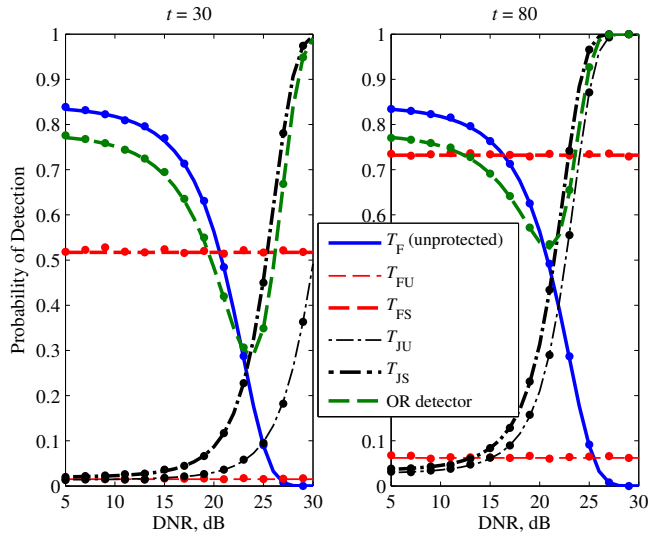


Fig. 1.  $P_{\text{D}}$  vs. DNR. Analytical (lines), simulation (markers).  $n = 400$ ,  $m = 20$ ,  $P_{\text{FA}} = 0.01$ ,  $\text{SNR} = 15$  dB.

taking  $P_{\text{FA},1} = P_{\text{FA},2} = 1 - \sqrt{1 - P_{\text{FA}}} \approx \frac{1}{2}P_{\text{FA}}$ , as this usually yields an acceptable tradeoff between the two individual tests. Regarding the probability of detection of the OR detector, even though independence of  $T_{\text{F}}$  and  $T_{\text{AS}}$  may not hold when under attack, we propose the approximation  $P_{\text{D}} \approx P_{\text{D},1} + P_{\text{D},2} - P_{\text{D},1}P_{\text{D},2}$ , with  $P_{\text{D},i}$  the detection probabilities of the individual tests.

## V. RESULTS

We consider a network with  $n = 400$  nodes and  $m = 20$ , setting  $P_{\text{FA}} = 0.01$ . It is assumed that  $\tilde{\mathbf{H}}$  is full rank ( $\tilde{r} = m$ ). Fig. 1 shows the detection probability of the different schemes for 15 dB SNR, as a function of the Distortion-to-Noise Ratio (DNR)  $\frac{\|\mathbf{a}\|^2}{\sigma^2}$  under orthogonal attacks  $\mathbf{a} \in \mathcal{R}^\perp(\mathbf{H})$ . Analytical results are shown together with those from Monte Carlo simulations (with entries of  $\mathbf{H}$  independently drawn from a  $\mathcal{N}(0, 1)$  distribution, and  $k = 22$  compromised sensors) for 30 and 80 trusted nodes. As expected, the unprotected detector (4), which does not account for attacks, performs best for low DNR but its detection probability goes to zero as DNR increases. The remaining schemes present different performance tradeoffs between the low- and high-DNR regimes. Since it only uses data from trusted sensors, the detector  $T_{\text{FU}}$  from (8) performs poorly if  $t/n$  is small. The behavior of  $T_{\text{JU}}$  from (16) and  $T_{\text{JS}}$  from (19) is better for high DNR, but not in the low DNR region for small  $t/n$ . In contrast, the detector  $T_{\text{FS}}$  from (11) and the OR detector from Sec. IV-C provide a much better performance-security tradeoff. Whereas  $T_{\text{FS}}$  is insensitive to orthogonal attacks, it may present a sizable loss for low DNR if  $t/n$  is small. Although the OR detector is affected by attacks, it performs well in the low- and high-DNR regions, and only for an interval of medium DNR values it is outperformed by the  $T_{\text{FS}}$  detector. The proposed analytical approximation for the probability of detection of the OR scheme is seen to accurately match the empirical results. This was the case in all settings investigated.

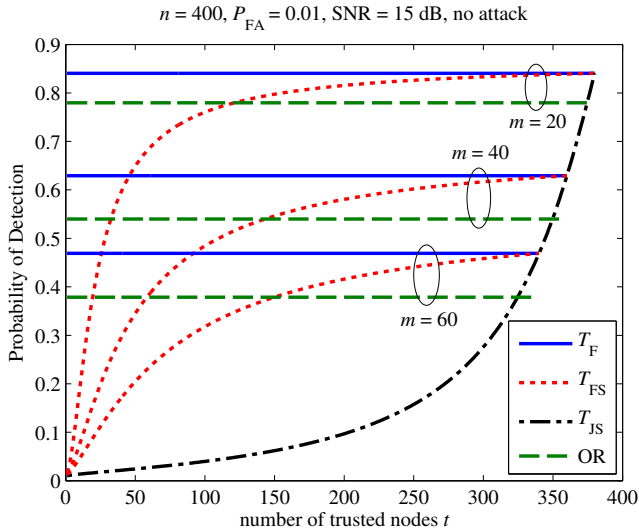


Fig. 2.  $P_D$  vs.  $t$  in the absence of attacks.

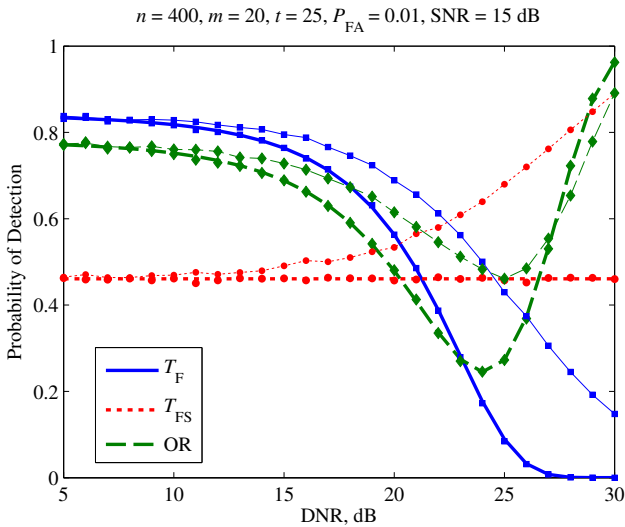


Fig. 3.  $P_D$  for orthogonal (thick) and blind attacks (thin lines).

To better gauge the improvement of the proposed schemes with the number of trusted nodes, Fig. 2 plots  $P_D$  vs.  $t$  at 15 dB SNR in the absence of attacks. This shows the loss with respect to the  $T_F$  benchmark incurred in the low DNR region (the price paid for robustness against attacks). In this regime, the loss of the OR detector is constant with  $t$  (for this detector, the benefit of a larger  $t$  is reflected in the medium DNR regime, in which the "dip" in the  $P_D$  vs. DNR curve becomes less pronounced as  $t$  increases, as seen in Fig. 1). The performance of  $T_{JS}$  improves very slowly with  $t$ . The improvement is significantly faster for  $T_{FS}$ , although the number of trusted nodes it requires to outperform the OR detector is relatively large (of the order of  $\frac{n}{3}$ ). Thus, the OR detector is a good choice when  $t/n$  is small; otherwise, the  $T_{FS}$  detector may be preferable.

We also checked (via Monte Carlo) the impact of *blind attacks*: entries of  $\mathbf{a}$  for compromised nodes are  $a_i \sim \mathcal{N}(0, 1)$  i.i.d., and scaled to meet a given DNR. As Fig. 3 shows,

even for detectors which assume orthogonal attacks ( $T_{FS}$  and OR), with blind attacks the probability of detection is actually larger than with orthogonal attacks (because the adversary cannot successfully exploit the attack component in the signal subspace), which justifies such worst-case assumption.

## VI. CONCLUSIONS

With trusted nodes available, robust detectors against malicious attacks have been developed and analyzed under the umbrella of matched subspace detection. Exploiting the structure of the attack vector may significantly reduce the number of required trusted nodes. The  $T_{FS}$  and OR detectors emerge as attractive choices with different performance-robustness tradeoffs. Our results can be used in network design to evaluate these tradeoffs, offering guidelines as to how many sensors must be secured for a given performance target.

## REFERENCES

- [1] P. Corke et al., "Environmental wireless sensor networks," *Proc. IEEE*, vol. 98, no. 11, pp. 1903–1916, Nov. 2010.
- [2] J.-F. Chamberland and V. V. Veeravalli, "Wireless sensors in distributed detection applications," *IEEE Signal Process. Mag.*, vol. 24, no. 3, pp. 16–25, May 2007.
- [3] A. Perrig, J. Stankovic, and D. Wagner, "Security in wireless sensor networks," *Commun. ACM*, vol. 47, no. 6, pp. 53–57, Jun. 2004.
- [4] Xiangqian Chen, K. Makki, Kang Yen, and N. Pissinou, "Sensor network security: a survey," *IEEE Commun. Surveys & Tutorials*, vol. 11, no. 2, pp. 52–73, 2nd quarter 2009.
- [5] M. C. Vuran, Ö. B. Akan, and I. F. Akyildiz, "Spatio-temporal correlation: theory and applications for wireless sensor networks," *Computer Networks*, vol. 45, pp. 245–259, 2004.
- [6] C. Guestrin et al., "Distributed regression: an efficient framework for modeling sensor network data," in *Int. Symp. Info. Process. Sensor Networks (IPSN)*, 2004, pp. 2394–2398.
- [7] A. Ribeiro and G. Giannakis, "Bandwidth-constrained distributed estimation for wireless sensor networks—Part II," *IEEE Trans. Signal Process.*, vol. 54, pp. 2784–2796, Jul. 2006.
- [8] A. Dogandžić and K. Qiu, "Decentralized random-field estimation for sensor networks using quantized spatially correlated data and fusion-center feedback," *IEEE Trans. Signal Process.*, vol. 56, no. 12, pp. 6069–6085, Dec. 2008.
- [9] R. López-Valcarce and D. Romero, "Design of data-injection adversarial attacks against spatial field detectors," in *IEEE Workshop Stat. Signal Process. (SSP)*, 2016.
- [10] S. Cui et al., "Coordinated data-injection attack and detection in the smart grid," *IEEE Signal Process. Mag.*, vol. 29, no. 5, pp. 106–115, Sep. 2012.
- [11] J. Kim, L. Tong, and R. J. Thomas, "Subspace methods for data attack on state estimation: a data driven approach," *IEEE Trans. Signal Process.*, vol. 63, pp. 1102–1114, Mar. 2015.
- [12] A. Pfitzmann, B. Pfitzmann, M. Schunter, and M. Waidner, "Trusting mobile user devices and security modules," *IEEE Computer*, vol. 30, no. 2, pp. 61–68, 1997.
- [13] P. Schaffer, K. Farkas, A. Horvath, T. Holczer, and L. Buttyán, "Secure and reliable clustering in wireless sensor networks: A critical survey," *Computer Networks*, vol. 56, no. 11, pp. 2726–2741, Jul. 2012.
- [14] R. López-Valcarce and D. Romero, "Attack detectors for data aggregation in clustered sensor networks," in *European Signal Process. Conf. (EUSIPCO)*, 2015.
- [15] S. Ganerwal, L. K. Balzano, and M. B. Srivastava, "Reputation-based framework for high integrity sensor networks," *ACM Trans. Sensor Networks*, vol. 4, no. 3, May 2008.
- [16] J. López, R. Román, I. Agudo, and C. Fernández-Gago, "Trust management systems for wireless sensor networks: Best practices," *Computer Communications*, vol. 33, no. 9, pp. 1086–1093, Jun. 2010.
- [17] L. L. Scharf and B. Friedlander, "Matched subspace detectors," *IEEE Trans. Signal Process.*, vol. 42, no. 8, pp. 2146–2157, Aug. 1994.
- [18] N. L. Johnson, S. Kotz, and N. Balakrishnan, *Continuous Univariate Distributions*, vol. 2, John Wiley, 2nd edition, 1995.