



**HAL**  
open science

# Distributed Bias Detection in Cyber-Physical Systems

Simon Thougard, Bruce Mcmillin

► **To cite this version:**

Simon Thougard, Bruce Mcmillin. Distributed Bias Detection in Cyber-Physical Systems. 14th International Conference on Critical Infrastructure Protection (ICCIP), Mar 2020, Arlington, VA, United States. pp.245-260, 10.1007/978-3-030-62840-6\_12 . hal-03794640

**HAL Id: hal-03794640**

**<https://inria.hal.science/hal-03794640v1>**

Submitted on 3 Oct 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



This document is the original author manuscript of a paper submitted to an IFIP conference proceedings or other IFIP publication by Springer Nature. As such, there may be some differences in the official published version of the paper. Such differences, if any, are usually due to reformatting during preparation for publication or minor corrections made by the author(s) during final proofreading of the publication manuscript.





## Chapter 1

# DISTRIBUTED BIAS DETECTION IN CYBER-PHYSICAL SYSTEMS

**Abstract** In distributed cyber-physical systems with noisy measurements, an attacker can effectively publish false measurements. These biased false measurements can be impossible to individually distinguish from noise and allow an attacker to gain a small but persistent economic advantage. This paper proposes to employ the residual sum, as a fundamental measurement of bias in cyber-physical systems, and develop a detection scheme for bias attacks. We discuss how the residual sum can both be highly efficient, and privacy preserving, while effectively detecting any bias attack.

tm

**Keywords:** Cyber-physical Systems, Security, Privacy, Smart-grid, False Data Injection, State Estimation

## 1. Introduction

False Data Injection Attacks (FDIA) on power systems have been the subject of intense study since introduced by [3]. The attack model assumes an attacker that knows the power system configuration, and has the ability to send corrupted measurements to a control entity, i.e. bad data injection. Furthermore, [3] proves that such attacks can be undetectable by standard methods.

FDIA's pose a fundamental challenge to any Cyber-Physical System (CPS): If a node in a CPS is compromised by an attacker, and the attacker knows what security measures are in place, the attacker can always inject bad data into a control system while avoiding detection. In this paper we will prove this for any CPS that is tolerant to measurement error. A more formal definition of FDIA, and their application to CPSs, will be given in section 3.

The effectiveness and limitations of such attacks was discussed and studied in the original paper; later papers have proposed defense schemes and variations of FDIA. Some proposed defence schemes suffer from a

simple lack of imagination. If the attacker *knows* of the defence schemes, can they still circumvent them? After all, in [3] the attacker is assumed to know both the system configuration and bad data thresholds.

This work approaches the problem in a more general sense. Sound defence schemes result from specific criteria. The proposed scheme confronts economic attacks specifically, and meets relevant criteria. To optimally balance the trade-off between false-positive rate and false-negative rate, a novel queue-based approach to attack detection is also proposed. Attacks on the smart grid specifically are considered. Conventional state estimation assumes a central system operator, who may be able to counteract an attack if properly identified. Under a distributed electric grid architecture, a centralized entity may still exist, but it is relevant to consider both privacy issues, as well as practical applicability of any defence scheme. The literature on FDIA represents the systems as matrices of data, but for any individual node, only a slice of that data will be available.

Section 2 contains a brief overview of existing work on FDIA. Section 3 outlines a formal definition of FDIA and CPSs and lays out problems with existing approaches along with arguments against some of the conventions in the literature. In section 4 contains a proposal for a simple yet effective approach to the specific attack goal of "theft". We derive bounds on the behavior of the solution. Section 5 contains simulation results.

## 2. Other Work

[3] introduced the concept of False Data Injection Attacks in power systems state estimation, and formalized a proof of existence of zero-residual attacks. The paper does not propose a solution, but derives expressions for an optimal attack vector under different conditions. While the zero-residual attack is the most impressive version of FDIA, it will not be considered in this paper. Zero-residual attacks can be considered unsolvable, as they result from an attacker with complete power to arbitrarily inject bad data. However, good system design may make such attacks very difficult to carry out. The attacker has to compromise every measurement related to a state variable. It can be compared to an attacker buying a bank, in order to get access to the vault. It is theoretically possible, but perhaps not a practical security concern. Therefore the focus is instead on attacks whose residual is below a tolerable threshold.

[4] gives a thorough review of the literature on FDIA in power systems. It may seem strange that so much work has been done on this subject, yet papers continue to be published on it. What remains to be studied? A lot of the work is split into different variations of the attack and the system under attack, as well as different approaches to defence schemes. However it seems preference is given to attack scenarios that are easy to define mathematically, rather than attack scenarios that capture rational behavior of an attacker. This stems no doubt from the academic norms of the control theory and power systems communities. As a result, an attacker that simply behaves in a sub-optimal or nonconforming way, may go perfectly undetected by detection schemes that are designed assuming optimal or conforming behavior.

[8] introduces economic attacks on electricity markets, using FDIA, and [2] proposes a solution to detect such attacks. Economic attacks differ from most of the literature by considering an attacker with a known and quantifiable goal. Such attacks may be designed to go undetected at the expense of some magnitude of the attack. If an entity gains a small but reliable amount of value from persistent subtle attacks, they would want to carry out the attack for as long as possible. However, it is not critical to quickly uncover such attacks, since they pose no security threat. It would be sufficient to guarantee eventual detection.

The Smart Grid presents many new opportunities and challenges, given the ability of nodes to coordinate production, consumption and distribution of electricity. [5] presents a comprehensive and decentralized approach for how to implement such coordination in a decentralized way. Privacy concerns of smart meters are considered in [6]. In the paper, smart meter data is used to infer private information about a home. The paper is a reminder of how secondary data sources, can reveal private sensitive information.

### 3. Problem Definition

To understand the context of FDIA, this paper will focus on the application of power systems, even though similar attacks can apply to any CPS, such as industrial control systems or positioning systems. In power systems control, a central operator collects measurements from the system, and decides whether or not to take some action. In a decentralized grid, such activity may be carried out in a distributed manner, where local nodes collaborate on control and decision making. Since the measurements contain some error, State Estimation (SE) is employed to calculate the most probable real state of the system. SE relies on the relations between states. This can be expressed through the model

$z = h(x) + e$ , where  $z$  is a measurement,  $x$  is a state in the system,  $h(x)$  is a function relating states to measurements and  $e$  is some error. SE is the problem of estimating the state  $\hat{x}$  from  $z$ , when there are multiple interrelated states. In this paper the details on how to calculate SE are ignored, and SE is expressed as a function of  $z$ :

$$SE(z) = \hat{x}$$

While not exactly trivial, the actual method of SE is inconsequential to the rest of the work in this paper. Suffice it to say, that given noisy measurements, some estimates of the state of the system can be produced. This is especially useful when considering a smart grid, in which coordination between nodes may be distributed, and traditional SE no longer applies.

Bad measurement detection, is the problem of determining if any measurement is abnormal, which could be simply an anomaly. It is calculated using the residual:

$$r = z - \hat{x}$$

which is simply the difference between measurement and estimate. Suppose the residual is tested against a threshold value  $r \leq \tau$ , this can determine if the measurement is believable or not.

### 3.1 Attack Model

A traditional FDIA attacker is assumed to have control over one or several measurements, and detailed knowledge of the layout of the system. The attack is modeled by  $z_a = h(x) + e + \alpha$ , where  $\alpha$  is some non-zero value. The assumption on the attacker is that an undetected value of  $\alpha$  will yield a gain, while  $-\alpha$  will yield a loss.

The goal of the attack is then in two parts: To avoid any detection scheme deployed by the controller and to effect some change in the estimated states.

Avoiding detection is a matter of keeping the residual  $r$  below some threshold. In [3], the existence of zero-residual attacks are proven. Those are attacks that effect change in  $\hat{x}$ , without changing  $r$ . It should be noted that such attacks are only possible if an attacker controls every measurement related to some state. Another type of attack is also explored, in which  $r$  may change, but kept under the threshold  $\tau$ . These have less impact on  $\hat{x}$ , but can be carried out with even a single corrupt measurement. The rest of this paper will only consider the latter type of attack, where an attack results in some change in the residual.



The goal of the attacker is then to maximize the change in  $\hat{x}$ , while satisfying  $r \leq \tau$ . Given that the attacker knows how  $SE(z)$  is calculated and how  $\tau$  is set, it is not very difficult to determine optimal  $\alpha$  to inject.

This work assumes an attacker who is attempting to inject a consistent, yet unobtrusive, bias. The bias may be relatively small, even compared to measurement variance. The attacker will try to "hide in the error" so to speak, by keeping the attack residual too small to be distinguishable from measurement error, but consistently in a direction that benefits the attacker. The benefit could be in the form of simple theft, or overcharging for the amount of supplied energy.

### 3.2 A Graph Approach

SE is traditionally considered as an optimization problem, and power systems modeled as systems of linear equations. The collected data is represented as a vector of measurements, whose relations are expressed by a matrix.

There are two reasons to pursue a graph-model of the system instead: We are not going to consider any optimization problems, and therefore we do not benefit from the standard form. Additionally, in a distributed system, such as the smart grid, control of the system may itself be distributed. Collecting all measurements for analysis may yield practical issues, such as long delays and response times. It may also be a privacy concern. A smart-grid can be modeled as a graph, where the nodes represent individual locations of the system, and the edges represent physical connections between states. Let  $G = V, E$  be a graph representation of the smart grid, where  $V$  is the set of nodes, representing an endpoint in the smart grid.  $E$  the set of edges, representing a transmission line between two nodes. Let  $N(v)$  denote set of neighbors of  $v$ .

By assuming that every node has an attached battery, able to store or release energy at will, the problem definition becomes simpler. To see that this is not a restriction, simply consider a node without a battery as a node that chooses to never use its battery. Each node will be able to measure and report its incoming power  $P_i^+$  from each neighbor, outgoing power  $P_j^-$  to each neighbor, and battery storage and discharge,  $P_b^+$  and  $P_b^-$  respectively. Each node will be under the constraint of conservation of power, expressed as:

$$\sum_{v_i \in N(v)} P_i^+ + \sum_{v_i \in N(v)} P_i^- + P_b^+ + P_b^- = 0$$

Each node is required to report incoming and outgoing power readings to each respective neighbor. To preserve privacy, we do not require that

these readings are reported to a central operator. Observe that even if a node is reporting all measurements to an observer, it is trivial for a node to seem perfectly consistent internally, even if it is reporting false external transmissions. If a node is under-reporting incoming power by  $a$ , it need simply subtract  $a$  from  $P_b^+$  to satisfy the equality.

To model economic attacks, where an attacker seeks to gain some advantage from an FDIA, it is useful to simplify the model further. Consider the case where two nodes,  $v_1$  and  $v_2$ , share a transmission line, and  $v_1$  intends to carry out a false data attack against  $v_2$ . There are two cases to consider: Either the attack is strictly directed toward  $v_2$ , in which case this attack can be considered a 2-node problem. Or the attack is directed against multiple nodes. In this case the attack will still produce some attack residual between  $v_1$  and  $v_2$ . In the case of faulty hardware, it may be sometimes reasonable to assume that a bad node provides bad data to all neighboring nodes. However, an intelligent attacker may decide to only provide bad data to a single neighbor, or a select set of neighbors. Hence it makes sense to reduce attacks to a 2-node problem. This may not be the optimal way to detect all attacks, as some attacks may be detected faster by considering multiple residuals. But setting optimality aside, it makes for a much simpler problem definition:

”Given a set of edges  $E$  of a cyber-physical system graph, determine which edges are likely to be under attack.”

Note the emphasis on edges, not nodes, which differs from the standard problem formulation, which emphasises nodes or state variables. This makes the problem explicitly about relations between nodes, not the nodes themselves. Residual sums represent relations between nodes, Attacks, when detected, reside in the relation between nodes.

### 3.3 Distinguishing Victim from Attacker

The problem formulation in the previous section only mentions determining an attacked edge, not which node may be the attacker. Although it may be possible in some circumstances to determine which node is the attacker, and which is the victim, it is impossible to do for all cases. Consider the case where an attacker at  $v_1$  carries out an attack strictly towards  $v_2$ . That is,  $v_2$  is the only node that can attest to the false data reported by  $v_1$ . In this case, a third party observer would only be able to conclude that a disagreement exists between  $v_1$  and  $v_2$ . In many practical situations we would require to know which node is acting falsely, but for the purpose of this paper, it is sufficient to determine the edge disagreement.

### 3.4 False Positives or False Negatives

A common approach to attack detection is to set a threshold for what data is considered normal, and what data is considered corrupt. This threshold can be statistically derived to have some desired property. Against a sophisticated attacker, such a threshold will also define which attacks are considered tolerable. If a threshold for divergence between two measurements is set to 10%, an attacker who controls one of the measurements may design the attack vector to approach the 10% divergence, without exceeding it. If this threshold is reduced to 1%, it would put a tighter limit on what attacks are tolerated, but it would also increase the rate false positives.

By adjusting the threshold for attack detection, the rate of false positives (FP) is adjusted. If the goal is fewer false positives, the threshold is increased, if the goal is to reduce the number of false negatives (FN), the threshold is lowered. The trade-off between the inversely correlated FP and FN is adjustable, but a practical solution ought to include an effective way to balance these.

### 3.5 Smart Attacker

If an attacker knows what detection schemes are employed, any solution must assume attackers that actively avoid detection. It is not sufficient to determine what the most optimal attack strategy is, and then defend against it. Sub-optimal attacks, would go completely unnoticed.

### 3.6 Smart Grid Example

To illustrate how this work applies to real world cyber-physical systems, consider the example of the smart grid. [7] lays out an example of how power mitigation works on distributed grid infrastructure. The individual nodes report how much power they consume and produce, and an observer then verifies that reported values are consistent. The model does not take noise into account, but we can imagine a simple approach to handle noise: If the difference between two related measurements exceeds a certain threshold, an observer will detect the anomaly and take appropriate action. If the reported measurements are within the distance threshold, the arithmetic mean is agreed upon as the true state.

The illustration in figure 1 shows how an attack may be carried out. Among two adjacent nodes, sharing a distribution line, it would be trivial to perform a bias attack for a malicious node. The attacker simply needs to under-report incoming power, when consuming, and over-report when

producing. The amount of bias depends on the threshold for tolerance, which the attacker is assumed to know. In effect, the attack will look like noise to any observer. The scenario is illustrated in figure 1, where the attacker under-reporting received power, which amounts to theft. The physical connection (1) shared between the nodes has an actual flow of power, that both nodes can measure (2). The nodes are expected to report the measurement to each other, but the attacker injects their bias (3). Both nodes calculate the mean state value (4), and residuals (5). If the residuals are small enough, no malicious activity is detected.

In the example from figure 1, the attacker has now effectively gained free power, by abusing noise tolerance. This works because nodes use consensus to determine the real state of the system. In the absence of attack, this would be the most reliable approach for many noisy distributed cyber-physical systems. The goal is therefore to demonstrate how bias attacks can be detected efficiently.

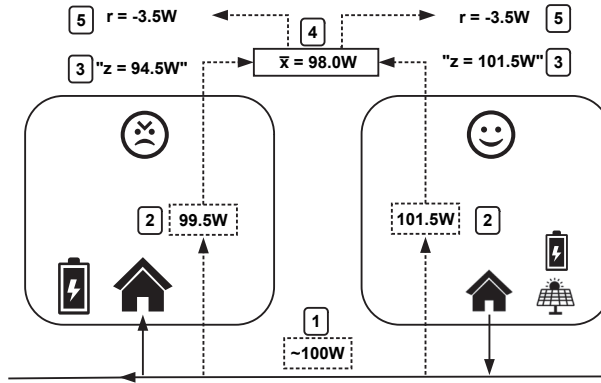


Figure 1. Illustration of bias attack on a smart grid power distribution line between two nodes

## 4. Proposed Solution

The following is an introduction and discussion of the residual sum, as a central measurement of the integrity of the system.

### 4.1 Residual Sum

Recall that the residual is defined as

$$r = z - \hat{x}$$

where  $\hat{x}$  is some estimate of the state  $x$ , and  $z$  is some measurement of  $x$ . A large  $|r|$  is obviously an indication of error or false data. However, if

we assume  $r \sim \mathcal{N}(0, \sigma^2)$  under no attack, at any time step,  $r$  is expected to be non-zero, and we expect to detect  $|r| > \tau$  at a steady rate, where  $\tau$  is some threshold for the residual. That is, we expect a certain rate of false positives. If an attacker injects  $a = \frac{1}{2}\sigma^2$ , the rate of positives will increase, but how is an observer to determine that an attack is happening? Since  $\sigma^2$  is not known, this is not a trivial problem.

In a single instance, it may be impossible to discriminate between an attack, and a random event. But if we keep track of the residual over multiple time instances, we can get a clearer picture. If we take the sum of absolute residuals until time  $n$ , we get a monotonically increasing function. But by summing the signed values of  $r$ , we get the function:

$$RSUM(n) = \sum_{i=0}^n r$$

The residual sum can be used, much the same way the residual is used, to determine if reported measurements are within expected bounds. A large  $|RSUM|$  indicates error, or false data, but  $RSUM$  has some properties that make it ideal for our stated problem.

Section 4.5 contains a practical application of the residual sum using a queue.

## 4.2 Statistical Behavior

The residual under no attack is assumed to be  $r \sim \mathcal{N}(0, \sigma^2)$ , and under biased attack,  $r \sim \mathcal{N}(c, \sigma^2)$  where  $c$  is some constant.

The residual sum, which is just the addition of residuals, has the distribution

$$RSUM(n) \sim \mathcal{N}(0, n \cdot \sigma^2)$$

This follows from the fact that the sum of two normally distributed variables, is a normally distributed variable with mean of the sum of means, and variance of the sum of variances. Summing up  $n$  residuals, yields the distribution in (5).

The distribution of  $RSUM$  under biased attack is

$$RSUM_a(n) \sim \mathcal{N}(n \cdot c, n \cdot \sigma^2)$$

So over time, the mean of the biased  $RSUM_a$  grows at a rate of  $c$ , whereas the mean of the unbiased  $RSUM$  is expected to stay at 0.

This leads to the important observation: Too much of a good thing, can be a bad thing

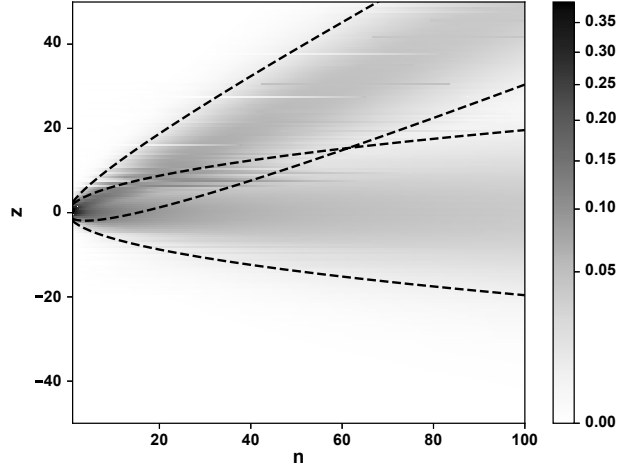


Figure 2. Diverging confidence intervals for biased and unbiased RSUM

The interpretation of this result, is that no matter what confidence level we choose ( $\gamma$ ), and how small the bias residual is, eventually the biased RSUM will distinguish itself from an unbiased RSUM. A small attack may go undetected for a long time, but eventually, it will be detectable with any arbitrary level of confidence.

### 4.3 Computation

*RSUM* is not only a useful tool for detecting bias attacks, but it is also requires very few resources to compute in a distributed system. In terms of space, note that it can be computed dynamically, and thus, each RSUM value requires only a single numerical value to store. The solution require each node to store  $2 \cdot \text{deg}(v)$ , where  $\text{deg}(v)$  is the degree of the node  $v$ . In terms of computation, each RSUM requires only one addition at each timestep. However, the estimate  $\hat{x}$  may require more computation, depending on method. In terms of bandwidth, each RSUM will require a single value transmitted between two nodes, at minimum. However, this value need only be transmitted between adjacent nodes, so it may not affect the communication network significantly.

### 4.4 Published Residual

One of the valuable properties of the residual sum, is the potential to protect privacy of individual nodes. In a traditional CPS, where data is sent to a central controller, privacy depends on how well the controller may be trusted. In a fully distributed system, this trust may not be

assumed, but with the residual sum it may not be an issue. Consider the case of the smart grid, where nodes may not wish to publish their energy consumption and production to a third party controller. If instead the residual sum is published, an observer would not be able to determine much. The residual sum only describes how much disagreement is between nodes, not actual amounts of energy transferred.

The only observation that can be made from the residual sum, is  $RSUM(n) < \sum_{i=0}^n(z)$ . This may be considered a privacy loss, since a non-zero RSUM would indicate activity, and vice versa. However, this concern can be alleviated by a noise-adding scheme:

Suppose two nodes share a physical connection, and publish their shared residual. To obfuscate activity between them, they add some noise value  $e$  to every published value.  $e$  is drawn from a list of values created in the following way:

- 1 A random seed value is exchanged between nodes
- 2 At coordinated intervals,  $n$  values are randomly generated by a Linear congruential generator, and added to a list  $l$
- 3 For every original element  $e$  in  $l$ , replace with  $\sigma \cdot (e\%100)/100$ , then add  $-e$  to the end of  $l$
- 4 For every  $2n$  timesteps, a random element  $e'$  is drawn and removed from  $l$ .  $e'$  is then added to the current published RSUM.

This scheme makes it impossible for an observer to determine if published values between time 0 and  $2n$  reflect activity or inactivity. Since the values of  $l$  sum up to 0, it also ensures that the published RSUM is accurate at the end of each interval. Although  $\sigma$  may not be a known value, some suitable scalar may be chosen in its place.

## 4.5 Action and Queue-based inspection

This subsection proposes how to actually use the residual sum to discover attacks. The solution assumes that some inspectors are tasked with finding and handling attacks and abnormalities. For electric grids, such inspectors already exist. They carry out routine inspections looking for meters that have been tampered with. It is assumed that inspectors are fully capable of both identifying and handling tampering at any location they visit, having full access to any part of the system and carrying specialized equipment. While this is a very convenient assumption, it is not far from reality. The main issue is resources and distribution, i.e. hiring inspectors and scheduling them effectively.

The proposal is to use the residual sum to prioritize the work of the inspectors. Simply put, focusing on those edges in the system that have a high residual sum. Since residual sums are published, the inspectors do not need to be employed locally, but can be employed by a central control agency for a large region.

One way to do it, would be to set a threshold,  $\tau$ , for residual sums, and then call inspectors to any edge that exceeds  $\tau$ . By choosing a  $\tau$  value, we would effectively be setting the rate of positives generated by the system. If we set  $\tau$  too low, we would generate more false positives than the inspectors could handle. If we set  $\tau$  too high, we would end up waiting longer than necessary, before acting on malfeasance. A good threshold would yield a manageable rate of positives, true and false.

However, a simpler approach exists, which does not rely on thresholds. Simply sort the residual sums by magnitude, and then schedule inspections at the corresponding locations in descending order. This way, the most likely attacked nodes are inspected first, and at the exact rate that the inspection crews can handle.

## 5. Results

To demonstrate the effectiveness of using the residual sum to detect bias attacks, we have conducted simulations. The simulations are carried out on a model of the IEEE 14-bus system found in [1], using the MATPOWER package for MatLab. To be able to artificially set the variance of the measurements, first a set of base measurements were produced by iteratively running state estimation on an initial set of measurements, and replacing them with the estimated values. This yields a measurement set with a square sum residual close to 0. These measurements serve as the basis for the simulations.

The simulations are run over 10000 time steps. At each step, some  $\mathcal{N}(0, \sigma^2)$  distributed noise is added to the base measurement of a node, with  $\sigma = 0.1$ . All simulations use 10 unbiased nodes and a varying number of biased nodes ( $k$ ). For the biased nodes, an additional 0.01 is added to each measurement, exactly 1/10th of the standard deviation of the noise. State estimation of the power system is carried out using MATPOWER at each time step, using the adjusted measurements. The residual of each variable is then calculated as the difference between measurement and estimate, and each residual added to its corresponding residual sum. At regular intervals ( $h$ ), the highest magnitude residual sum is identified and reset to 0. These results of one of these are displayed in Figure 3, with  $k = 2$  and  $h = 500$ . In this particular simulation instance, 13 out of the 19 inspections found a biased node to have the



highest RSUM, with a true positive rate of 68%. The figure illustrates the chaotic nature of the noisy measurements, where some of the unbiased RSUMs do become outliers, while the biased nodes consistently grow in the positive direction.

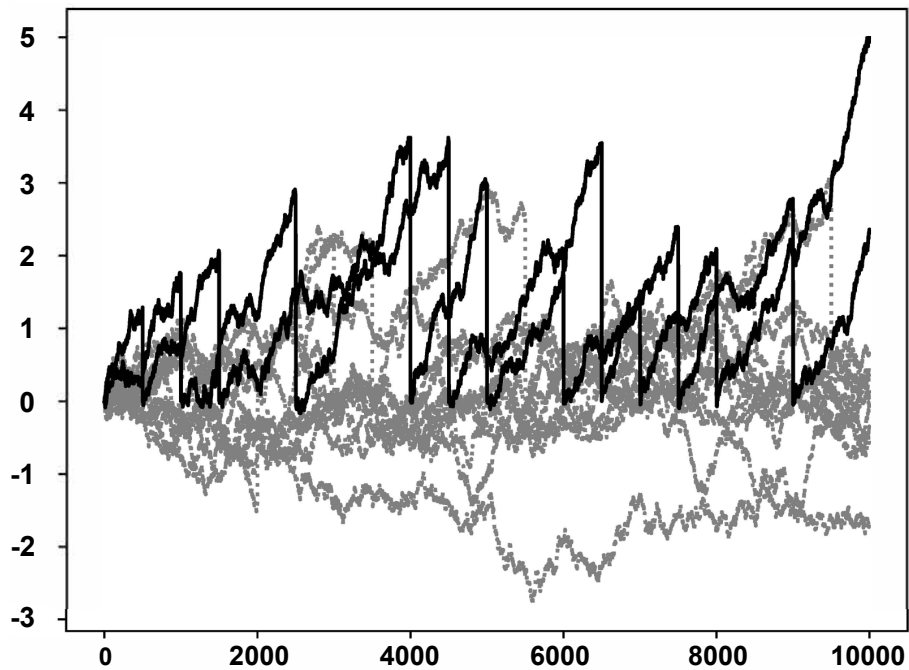


Figure 3. Simulation results with 10 unbiased RSUMs (gray), and 2 biased RSUMs (black). At regular intervals of 500 time steps, the highest residual sum is reset to 0.

Table 1 contains the true positive rates resulting from simulations with varying number of biased nodes and inspection intervals. Note that with long intervals and many biased nodes, the true positive rate effectively becomes 100%. This table illustrates one of the practical benefits of a queue approach. Inspectors will know the true positive rate and be able to adjust the interval, but will not know anything about biased nodes. By adjusting the inspection interval, they may achieve a desirable true positive rate.

This demonstrates the practical applicability of Theorem 1. The biased and unbiased nodes diverge, making it easy to prioritize which edges to focus inspection on.

k	Inspection Interval				
	100	300	500	700	900
1	0.27	0.45	0.47	0.57	0.73
2	0.42	0.64	0.68	0.64	0.91
3	0.55	0.67	0.74	0.86	0.91
4	0.7	0.82	0.95	1.0	1.0
5	0.74	0.88	0.95	1.0	1.0
6	0.78	0.91	1.0	1.0	1.0
7	0.86	0.94	1.0	1.0	1.0
8	0.86	1.0	1.0	1.0	1.0
9	0.83	0.94	1.0	1.0	1.0
10	0.93	1.0	1.0	1.0	1.0

Table 1. Rate of true positives with varying number of biased nodes (k) and inspection interval

## 6. Conclusion

A sophisticated attacker can easily carry out bias attacks on any noisy cyber-physical system, avoiding conventional detection methods. To detect such attacks, this paper proposes the residual sum as the basis for an approach, that specifically targets bias attacks. We have discussed the properties that make the residual sum optimal for this purpose, and derived some theoretical bounds on the behavior, both in the presence of attacks and under no attacks. Simulations are carried out, that demonstrate and visualize these assertions. It would also be interesting to carry out experiments in a physical testbed.

The specific application of electricity theft is presented, but future work may focus on other economic attacks. Additionally, multiple colluding attackers may be considered, which may lead to more challenges and opportunities.

Simon Thougard and Bruce McMillin This work was sponsored by a grant from the US National Science Foundation under award number CNS-1837472 and with support from the Missouri S&T Intelligent Systems Center.

## References

- [1] Illinois Center for a Smarter Electric Grid, IEEE 14-Bus System, ([icseg.itl.illinois.edu/ieee-14-bus-system/](http://icseg.itl.illinois.edu/ieee-14-bus-system/)).
- [2] L. Jia, R. Thomas and L. Tong, Malicious data attack on real-time electricity market, *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 5952–5955,

2011.

- [3] Y. Liu, P. Ning and M. Reiter, False data injection attacks against state estimation in electric power grids, *ACM Transactions on Information and System Security*, vol. 14(1), 2011.
- [4] G. Liang, J. Zhao, F. Luo, S. Weller and Z. Dong, A review of false data injection attacks against modern power systems, *IEEE Transactions on Smart Grid*, vol. 8(4), pp. 1630–1638, 2016.
- [5] E. Mengelkamp, B. Notheisen, C. Beer, D. Dauer and C. Weinhardt, A blockchain-based smart grid: Towards sustainable local energy markets, *Computer Science - Research and Development*, vol. 33(1-2), pp. 207–214, 2018.
- [6] A. Molina-Markham, P. Shenoy, K. Fu, E. Cecchet and D. Irwin, Private memoirs of a smart meter, *Proceedings of the Second ACM Workshop on Embedded Sensing Systems for Energy-Efficiency in Building*, pp. 61–66, 2010.
- [7] T. Roth and B. McMillin, Physical attestation of cyber processes in the smart grid, in *Critical Information Infrastructures Security*, E. Luijff and P. Hartel (Eds.), Springer, Cham, Switzerland, pp. 96–107, 2013.
- [8] L. Xie, Y. Mo and B. Sinopoli, False data injection attacks in electricity markets, *Proceedings of the First IEEE International Conference on Smart Grid Communications*, pp. 226–231, 2010.