

False-Data Attacks in Stochastic Estimation Problems with Only Partial Prior Model Information

Adrian N. Bishop

Abstract—The security of state estimation in critical networked infrastructure such as the transportation and electricity (smart grid) networks is an increasingly important topic. Here, the problem of recursive estimation and model validation for linear discrete-time systems with partial prior information is examined. Further, detection of false-data attacks on robust recursive estimators of this type is considered. The framework considered in this work is stochastic. An underlying linear discrete-time system is considered where the statistics of the driving noise is assumed to be known only partially. A set-valued estimator is then derived and the conditional expectation is shown to belong to an ellipsoidal set consistent with the measurements and the underlying noise description. When the underlying noise is consistent with the underlying partial model and a sequence of realized measurements is given then the ellipsoidal, set-valued, estimate is computable using a Kalman filter-type algorithm. A group of attacking entities is then introduced with the goal of compromising the integrity of the state estimator by hijacking the sensor and distorting its output. It is shown that in order for the attack to go undetected, the distorted measurements need to be carefully designed.

I. INTRODUCTION

The Kalman filter is an optimal filter, in the minimum variance sense, in the class of linear filters; see [1]. Much of the theory on Kalman filtering presumes the noise driving the process and the measurement sequence is Gaussian and white. Of course, this is not a necessary prerequisite for optimality. However, if one does assume Gaussian noise processes then it is of interest to know what the effect of applying incorrect noise statistics are on the performance of the filter. For example, early work in this area was proposed by [2] where the effect of incorrect noise covariances are used in the Kalman filtering algorithms and the performance is analysed. Related work and analysis is given in [3], [4], [5], [6]. This list is by no means exhaustive.

The problem of state estimation with uncertainty in the model has been widely investigated in the field of robust control and filtering. In this field it is typical to assume the system and measurement model itself has uncertainty and to model this uncertainty as a noise input drawn from a particular class (or set) of signals. Early work along these lines is given in [7], [8], [9], [10] where min-max type and set-valued type filtering results are related to the Kalman filter. These papers were generalized in [11], [12], [13], [14], [15] where the uncertainty was characterised by integral quadratic constraints. Similar work combining stochastic and set-based uncertainties has also been considered; e.g. see [16]. As discussed in, e.g., [17], set-valued state estimation is particularly suited to a number of applications such as target tracking. There are often physical constraints on the set of

target states, e.g. the target speed might be upper bounded etc, which leads to non-Gaussian target state distributions; see [17].

Bayesian algorithms for posterior estimation and partial prior knowledge etc have also been considered. Again, the use of set-based methods have been used; see [18], [19]. The work of [20] is an underlying motivation for much of the Bayesian development. Later, we outline a different method for introducing stochastic uncertainties and dealing with partial models of the statistics.

The idea of set-based estimation and robust filtering has a natural relationship with the model-validation problem. More specifically, for classes of uncertainties, set-valued state estimators can be used to determine if the measured data is consistent with the assumed system model; e.g. by checking if the estimator (or some statistic) falls within the derived set. This form of model validation was explored in, e.g., [11], [12], [13], [14], [15] and is further explored in this paper where stochastic uncertainties are considered.

The first high-level problem considered in this work is that of robust recursive Bayesian filtering given partial prior information about the noise statistics. This partial prior information comes in the form of a nominal system model and a particular kind of constraint. In particular, the deviation of the actual, real-world, model being observed is modelled by a constraint on a probability measure obtained via a particular change-of-measure operation on the nominal system's measure.

Additionally, this work is concerned with the problem of *securing estimation and control systems*. More specifically, this paper follows [21], [22], [23], [24] and considers the problem of safeguarding state estimators in critical infrastructure. Specifically, we study false-data attacks on stochastic and robust state estimation. We will consider an underlying class of stochastically uncertain (discrete-time) systems and we will outline a set-valued state estimation algorithm that recursively produces an ellipsoidal set of all those expected state estimates consistent with the measurements and modelling assumptions. We then draw on this set to infer the probability of an attack on the system and to determine when one cannot safely produce a consistent expected estimate.

II. PROBLEM SETUP

The scenario introduced in this section is novel in the sense that we consider a robust estimation scenario involving stochastic uncertain systems formulated using Bayesian probability and change-of-measure theory.

A. The Nominal Model

Fix a probability space $(\mathcal{S}, \mathcal{F}, \mathbb{P})$. Consider a discrete-time stochastic system of the form

$$\mathbf{X}_{t+1} = \Phi_t \mathbf{X}_t + \mathbf{E}_{t+1} \quad (1)$$

$$\mathbf{Y}_t = \Gamma_t \mathbf{X}_t + \mathbf{N}_t \quad (2)$$

where $\mathbf{X} \in \mathbb{R}^n$ denotes the *state* of the system and $\mathbf{Y} \in \mathbb{R}^l$ is the *measured output* of the system. The notation $\mathbf{X} \in \mathbb{R}^n$ etc implies $\mathbf{X} : \mathcal{S} \rightarrow \mathbb{R}^n$. The *random noise inputs* to the system are $\mathbf{E} \in \mathbb{R}^n$ and $\mathbf{N} \in \mathbb{R}^l$. Also, a known *control input* $\mathbf{u} \in \mathbb{R}^m$ could be applied but is omitted for brevity.

NOTE: Throughout this article when ‘=’ is used with expressions of random variable (defined on the same probability space) on both sides it is used to mean that the random variable expression on the left side is *equal in distribution* to the random variable expression on the right side. When ‘ \triangleq ’ is used it means that we are defining a symbol (on the left side) to represent the collection of symbols or expression appearing on the right side.

In many cases we could legitimately think of ‘=’ in a stronger sense but such a strengthening is typically unnecessary (unless otherwise noted).

Assumption 1. Let $T > 0$ be some terminal time, $t \in [0, T]$, and let

$$\mathfrak{G} \triangleq [\mathbf{X}_0^\top \ \mathbf{E}_1^\top \ \dots \ \mathbf{E}_t^\top \ \mathbf{N}_0^\top \ \mathbf{N}_t^\top]^\top \left(\{d\mathbf{x}_0 \times d\mathbf{e}_1 \times \dots \times d\mathbf{e}_t \times d\mathbf{n}_0 \times \dots \times d\mathbf{n}_t\} \right)$$

We have

$$\begin{aligned} \mathbb{P}(\mathfrak{G}) &= \frac{d\mathbf{x}_0}{\sqrt{(2\pi)^n |\Xi|}} e^{-\frac{1}{2}(\mathbf{x}_0 - \bar{\mathbf{x}})^\top \Xi^{-1}(\mathbf{x}_0 - \bar{\mathbf{x}})} \times \\ &\prod_{k=1}^t \frac{d\mathbf{e}_k}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(\mathbf{e}_k)^\top \Sigma^{-1}(\mathbf{e}_k)} \times \\ &\prod_{k=0}^t \frac{d\mathbf{n}_k}{\sqrt{(2\pi)^l |\Omega|}} e^{-\frac{1}{2}(\mathbf{n}_k)^\top \Omega^{-1}(\mathbf{n}_k)} \end{aligned} \quad (4)$$

Let $\mathcal{F}_t \subseteq \mathcal{F}_{t+1} \subseteq \mathcal{F}$ denote the completed σ -algebra generated by

$$\sigma(\{\mathbf{Y}_0, \mathbf{X}_0, \dots, \mathbf{Y}_t, \mathbf{X}_t\}) \quad (5)$$

and let \mathcal{Y}_t denote the completed σ -algebra generated by

$$\sigma(\{\mathbf{Y}_0, \dots, \mathbf{Y}_t\}) \quad (6)$$

Define a so-called innovation sequence $\{\tilde{\mathbf{N}}_t\}$ by

$$\tilde{\mathbf{N}}_t \triangleq \mathbf{Y}_t - \Gamma_t \mathbb{E}_{\mathbb{P}}[\mathbf{X}_t | \mathcal{Y}_{t-1}] \quad (7)$$

and let $\tilde{\mathcal{N}}_t$ denote $\sigma(\{\tilde{\mathbf{N}}_0, \dots, \tilde{\mathbf{N}}_t\})$.

Lemma 1. $\mathcal{Y}_t = \tilde{\mathcal{N}}_t$.

The proof of this lemma has been considered previously, e.g. see [25], [26], and similar results are given later in this paper without proof.

Note that $\tilde{\mathbf{N}}_t$ has zero mean and covariance matrix

$$\Gamma_t \text{Cov}_{\mathbb{P}}(\mathbf{X}_t - \mathbb{E}_{\mathbb{P}}[\mathbf{X}_t | \mathcal{Y}_{t-1}]) \Gamma_t^\top + \Omega \quad (8)$$

where $\text{Cov}_{\mathbb{P}}(\mathbf{X}_t - \mathbb{E}_{\mathbb{P}}[\mathbf{X}_t | \mathcal{Y}_{t-1}])$ denotes the covariance of $\mathbf{X}_t - \mathbb{E}_{\mathbb{P}}[\mathbf{X}_t | \mathcal{Y}_{t-1}]$ under \mathbb{P} .

B. A Change of Measure

The relative entropy of two probability measures \mathbb{P} and \mathbb{Q} is defined by

$$h(\mathbb{Q} || \mathbb{P}) \triangleq \begin{cases} \mathbb{E}_{\mathbb{Q}} \left[\log \left(\frac{d\mathbb{Q}}{d\mathbb{P}} \right) \right] & \text{if } \mathbb{Q} \ll \mathbb{P} \text{ and} \\ \log \left(\frac{d\mathbb{Q}}{d\mathbb{P}} \right) \in L^1(\mathbb{Q}) & \\ \infty & \text{otherwise} \end{cases} \quad (9)$$

where $\mathbb{E}_{\mathbb{Q}}[\cdot]$ denotes the expectation with respect to \mathbb{Q} and $\mathbb{Q} \ll \mathbb{P}$ means \mathbb{Q} is absolutely continuous with respect to \mathbb{P} .

Suppose that there exists two \mathcal{F}_t -adapted random sequences $\{\zeta_t\}$ and $\{\eta_t\}$. Then, suppose that the restriction of \mathbb{Q} on \mathcal{F}_t satisfies

$$\begin{aligned} \frac{d\mathbb{Q}}{d\mathbb{P}} \Big|_{\mathcal{F}_t} &= \Psi_t = e^{-\frac{1}{2} \mathbf{a}^\top \Xi^{-1} \mathbf{a} + \mathbf{a}^\top \Xi^{-1} (\mathbf{x}_0 - \bar{\mathbf{x}})} \times \\ &\prod_{k=1}^t e^{-\frac{1}{2} (\zeta_k)^\top \Sigma^{-1} (\zeta_k) + (\zeta_k)^\top \Sigma^{-1} (\mathbf{e}_k)} \times \\ &\prod_{k=0}^t e^{-\frac{1}{2} (\eta_k)^\top \Omega^{-1} (\eta_k) + (\eta_k)^\top \Omega^{-1} (\mathbf{n}_k)} \end{aligned} \quad (10)$$

(3)

where $\mathbf{a} \in \mathbb{R}^n$ and $t \in [0, T]$. Note that Ψ_t is a strictly-positive \mathcal{F}_t -measurable random variable and $\mathbb{E}_{\mathbb{P}}[\Psi_t] = 1$. Now it follows that

$$\begin{aligned} \mathbb{Q}(\mathfrak{G}) \Big|_{\mathcal{F}_t} &= \frac{d\mathbf{x}_0}{\sqrt{(2\pi)^n |\Xi|}} e^{-\frac{1}{2} (\mathbf{x}_0 - \bar{\mathbf{x}} - \mathbf{a})^\top \Xi^{-1} (\mathbf{x}_0 - \bar{\mathbf{x}} - \mathbf{a})} \times \\ &\prod_{k=1}^t \frac{d\mathbf{e}_k}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2} (\mathbf{e}_k - \zeta_k)^\top \Sigma^{-1} (\mathbf{e}_k - \zeta_k)} \times \\ &\prod_{k=0}^t \frac{d\mathbf{n}_k}{\sqrt{(2\pi)^l |\Omega|}} e^{-\frac{1}{2} (\mathbf{n}_k - \eta_k)^\top \Omega^{-1} (\mathbf{n}_k - \eta_k)} \end{aligned} \quad (11)$$

and the sequences $\{\mathbf{U}_t\} = \{\mathbf{E}_t - \zeta_t\}$ and $\{\mathbf{V}_t\} = \{\mathbf{N}_t - \eta_t\}$ on $(\mathcal{S}, \mathcal{F}, \mathbb{Q})$ are Gaussian with zero-mean and variances Σ and Ω respectively.

Now

$$\begin{aligned} h(\mathbb{Q} || \mathbb{P}) &= \frac{1}{2} \mathbf{a}^\top \Xi^{-1} \mathbf{a} + \frac{1}{2} \sum_{t=1}^T \mathbb{E}_{\mathbb{Q}} \left[(\zeta_t)^\top \Sigma^{-1} (\zeta_t) \right] \\ &\quad + \frac{1}{2} \sum_{t=0}^T \mathbb{E}_{\mathbb{Q}} \left[(\eta_t)^\top \Omega^{-1} (\eta_t) \right] \end{aligned} \quad (12)$$

is the relative entropy introduced into the system as a result of the measure change on the interval $t \in [0, T]$.

Definition 1 (Simple Energy Constraint). Suppose $\Xi = \Xi^\top > 0$, $\Sigma = \Sigma^\top > 0$, $\Omega = \Omega^\top > 0$, $\mathbf{a} \in \mathbb{R}^n$ and $\delta > 0$. Also there exists a finite time interval $[0, T]$. Then,

we consider the class of uncertain inputs $\{\zeta, \eta\}$ and initial conditions \mathbf{a} such that

$$2h(\mathbf{Q}||\mathbf{P}) \leq \delta \quad (13)$$

holds. Let \mathcal{Q}_{sec} denote the set of measures \mathbf{Q} such that $h(\mathbf{Q}||\mathbf{P}) < \infty$ and (13) holds.

In Definition 1 the symbols $\zeta \in \mathbb{R}^n$ and $\eta \in \mathbb{R}^l$ are random variables.

A less rigorous, but practically similar, condition can be stated in terms of probabilistically bounded system disturbances, and such an approach has been used in a number of robust estimation problems; see [27], [28].

Before proceeding, for further insight, imagine a normally distributed random vector \mathbf{X} on $(\mathcal{S}, \mathcal{F}, \mathbf{P})$ with zero-mean and constant covariance \mathbf{A} . Now imagine another random vector $(\mathbf{X} + \mu)$ on $(\mathcal{S}, \mathcal{F}, \mathbf{P})$ where μ is constant. One can think of μ as acting on the set of outcomes. However, now let

$$d\mathbf{Q} = e^{-\frac{1}{2}\mu^\top \mathbf{A}^{-1}\mu + \mu^\top \mathbf{A}^{-1}\mathbf{x}} d\mathbf{P} \quad (14)$$

such that on $(\mathcal{S}, \mathcal{F}, \mathbf{Q})$ the random vector \mathbf{X} has mean μ and constant covariance \mathbf{A} . By changing the measure from \mathbf{P} to \mathbf{Q} in this case we are leaving the outcomes alone but we are assigning different probabilities to the sets of events in \mathcal{F} (and consequently to their outcomes).

C. A Stochastically Uncertain System Model

Consequently, the system (1) and (2) on the probability space $(\mathcal{S}, \mathcal{F}, \mathbf{Q})$ is

$$\mathbf{X}_{t+1} = \Phi_t \mathbf{X}_t + \mathbf{U}_{t+1} + \zeta_{t+1} \quad (15)$$

$$\mathbf{Y}_t = \Gamma_t \mathbf{X}_t + \mathbf{V}_t + \eta_t \quad (16)$$

and we note that the measure we are working under should be clear from the context and notation. Note that we have not specified the distribution of the \mathcal{F}_t -adapted random sequences $\{\zeta_t\}$ and $\{\eta_t\}$. We have simply, by constraining the set of admissible measure changes to \mathcal{Q}_{sec} , bounded their first moment in some fashion.

Lemma 2. Fix the probability space $(\mathcal{S}, \mathcal{F}, \mathbf{Q})$ and consider the system described by (15) and (16) with $\mathbf{Q} \in \mathcal{Q}_{sec}$. Then

$$\mathbb{E}_{\mathbf{Q}}[\mathbf{U}_t] = \mathbb{E}_{\mathbf{Q}}[\mathbf{V}_t] = 0 \quad (17)$$

and the variance under \mathbf{Q} of \mathbf{U}_t and \mathbf{V}_t is Σ and Ω respectively. Furthermore, $\mathbb{E}_{\mathbf{Q}}[\mathbf{X}_0] = -\mathbf{a}$ and the variance under \mathbf{Q} of \mathbf{X}_0 is Ξ .

This lemma is really a direct consequence of (10) and the definition of \mathbf{Q} . Note that the system (15) and (16) with $\mathbf{Q} \in \mathcal{Q}_{sec}$ corresponds to the *real world* system while the system (1) and (2) under \mathbf{P} corresponds to a *nominal system model*. This corresponds to the physical scenario in which we do not know the exact characteristics of the real world system but rather where we have a nominal system model and a model class of uncertainty that specifies how the real world model may deviate from the nominal model.

Proposition 1. The completed filtration \mathcal{F}_t on the system (1) and (2) under \mathbf{P} is equivalent to the completed filtration \mathcal{F}_t on (15) and (16) under \mathbf{Q} . Similarly, the completed filtration \mathcal{Y}_t on the system (1) and (2) under \mathbf{P} is equivalent to the completed filtration \mathcal{Y}_t on (15) and (16) under \mathbf{Q} .

This proposition simply states that no more or no less events are possible under \mathbf{Q} than under \mathbf{P} . Changing the measure simply changes the probability we assign to each event.

III. A BAYESIAN SOLUTION TO THE ROBUST FILTERING PROBLEM

The following standing assumption is adopted throughout this section.

Assumption 2. Both $\{\zeta_t\}$ and $\{\eta_t\}$ are sequences of degenerate random variables with $\mathbb{Q}(\zeta_t = \tilde{\zeta}_t) = 1$ and $\mathbb{Q}(\eta_t = \tilde{\eta}_t) = 1$ almost surely for constant values $\tilde{\zeta}_t \in \mathbb{R}^n$ and $\tilde{\eta}_t \in \mathbb{R}^l$.

Note that this assumption, while quite strong, does not invalidate the spirit of the problem. Now note that

$$h(\mathbf{Q}||\mathbf{P}) = \frac{1}{2} \mathbf{a}^\top \Xi^{-1} \mathbf{a} + \frac{1}{2} \sum_{t=1}^T \mathbb{E}_{\mathbf{Q}}[\zeta_t^\top \Sigma^{-1} \zeta_t] + \frac{1}{2} \sum_{t=0}^T \mathbb{E}_{\mathbf{Q}}[\eta_t^\top \Omega^{-1} \eta_t] \quad (18)$$

Under \mathbf{P} recall that $\tilde{\mathbf{N}}_t \triangleq \mathbf{Y}_t - \Gamma_t \mathbb{E}_{\mathbf{P}}[\mathbf{X}_t | \mathcal{Y}_{t-1}]$ and $\tilde{\mathcal{N}}_t = \mathcal{Y}_t$ from Lemma 1 and $\tilde{\mathbf{N}}_t$ is zero-mean with variance $\Gamma_t \mathbf{M}_{t|t-1} \Gamma_t^\top + \Omega$. Under \mathbf{Q} define another so-called innovation sequence $\{\check{\mathbf{N}}_t\}$ by

$$\check{\mathbf{N}}_t \triangleq \mathbf{Y}_t - \Gamma_t \mathbb{E}_{\mathbf{Q}}[\mathbf{X}_t | \mathcal{Y}_{t-1}] \quad (19)$$

and let $\check{\mathcal{N}}_t$ denote $\sigma(\{\check{\mathbf{N}}_0, \dots, \check{\mathbf{N}}_t\})$.

Lemma 3. $\mathcal{Y}_t = \check{\mathcal{N}}_t = \check{\mathcal{N}}_t$.

Under \mathbf{Q} it follows that $\check{\mathbf{N}}_t$ has mean $\mathbb{E}_{\mathbf{Q}}[\eta_t]$ and covariance matrix

$$\Gamma_t \text{Cov}_{\mathbf{P}}(\mathbf{X}_t - \mathbb{E}_{\mathbf{P}}[\mathbf{X}_t | \mathcal{Y}_{t-1}]) \Gamma_t^\top + \Omega \quad (20)$$

where

$$\text{Cov}_{\mathbf{P}}(\mathbf{X}_t - \mathbb{E}_{\mathbf{P}}[\mathbf{X}_t | \mathcal{Y}_{t-1}]) = \text{Cov}_{\mathbf{Q}}(\mathbf{X}_t - \mathbb{E}_{\mathbf{Q}}[\mathbf{X}_t | \mathcal{Y}_{t-1}]) \quad (21)$$

denotes the covariance of $\mathbf{X}_t - \mathbb{E}_{\mathbf{P}}[\mathbf{X}_t | \mathcal{Y}_{t-1}]$ under \mathbf{P} . The fact that the variance of $\check{\mathbf{N}}_t$ under \mathbf{Q} is the same as the variance of $\tilde{\mathbf{N}}_t$ under \mathbf{P} is easily verified.

We then have the following theorem.

Theorem 1. Let $t \in [0, T]$. The expectation $\mathbb{E}_{\mathbf{Q}}[\mathbf{X}_t | \mathcal{Y}_t]$ of (15) and (16) with $\mathbf{Q} \in \mathcal{Q}_{sec}$ belongs to an ellipsoidal set

$$\mathfrak{X}_t = \{ \mathbb{E}_{\mathbf{Q}}[\mathbf{X}_t | \mathcal{Y}_t] \in \mathbb{R}^n : \Upsilon_t^\top \mathbf{A}_{t|t}^{-1} \Upsilon_t \leq \delta - \kappa_t \} \quad (22)$$

where $\mathbf{A}_{t|t} \triangleq \text{Cov}_P(\mathbf{X}_t - \mathbb{E}_P[\mathbf{X}_t|\mathcal{Y}_t]|\mathcal{Y}_t)$ denotes the covariance of $\mathbf{X}_t - \mathbb{E}_P[\mathbf{X}_t|\mathcal{Y}_t]$ under P conditioned on \mathcal{Y}_t and where

$$\Upsilon_t \triangleq \mathbb{E}_Q[\mathbf{X}_t|\mathcal{Y}_t] - \mathbb{E}_P[\mathbf{X}_t|\mathcal{Y}_t] \quad (23)$$

and

$$\kappa_t \triangleq \sum_{k=0}^t \tilde{\mathbf{N}}_k^\top [\boldsymbol{\Omega} + \boldsymbol{\Gamma}_k \mathbf{A}_{t|t-1} \boldsymbol{\Gamma}_k^\top]^{-1} \tilde{\mathbf{N}}_k \quad (24)$$

where $\mathbf{A}_{t|t-1} \triangleq \text{Cov}_P(\mathbf{X}_t - \mathbb{E}_P[\mathbf{X}_t|\mathcal{Y}_{t-1}]|\mathcal{Y}_{t-1})$ denotes the covariance of $\mathbf{X}_t - \mathbb{E}_P[\mathbf{X}_t|\mathcal{Y}_{t-1}]$ under P conditioned on \mathcal{Y}_{t-1} and

$$\tilde{\mathbf{N}}_k \triangleq \mathbf{Y}_k - \boldsymbol{\Gamma}_k \mathbb{E}_P[\mathbf{X}_k|\mathcal{Y}_{k-1}] \quad (25)$$

and the centroid of the ellipse is given by $\mathbb{E}_P[\mathbf{X}_t|\mathcal{Y}_t]$.

Let $\mathbf{Y}_t = \{\mathbf{y}_0, \dots, \mathbf{y}_t\}$ with $\mathbf{y}_t \in \mathbb{R}^l$ denote the set of realized measurements. Consider the following Riccati equations

$$\mathbf{M}_{t|t} = \left[\mathbf{M}_{t|t-1}^{-1} + \boldsymbol{\Gamma}_t^\top \boldsymbol{\Omega}^{-1} \boldsymbol{\Gamma}_t \right]^{-1} \quad (26)$$

$$\mathbf{M}_{t|t-1} = \boldsymbol{\Phi}_{t-1} \mathbf{M}_{t-1|t-1} \boldsymbol{\Phi}_{t-1}^\top + \boldsymbol{\Sigma} \quad (27)$$

$$\mathbf{M}_{0|0} = \boldsymbol{\Xi} \quad (28)$$

where the existence of a positive-definite solution to $\mathbf{M}_{t|t}$ is guaranteed; see e.g. [29]. Consider also the following set of state equations

$$\hat{\mathbf{x}}_{t|t} = \boldsymbol{\Phi}_{t-1} \hat{\mathbf{x}}_{t-1|t-1} + \mathbf{M}_{t|t} \boldsymbol{\Gamma}_t^\top \boldsymbol{\Omega}^{-1} (\mathbf{y}_t - \boldsymbol{\Gamma}_t \boldsymbol{\Phi}_{t-1} \hat{\mathbf{x}}_{t-1|t-1}) \quad (29)$$

$$\hat{\mathbf{x}}_{0|0} = \bar{\mathbf{x}} \quad (30)$$

and

$$\hat{\kappa}_t = \hat{\kappa}_{t-1} + \nu_t^\top [\boldsymbol{\Gamma}_t \mathbf{M}_{t|t-1} \boldsymbol{\Gamma}_t^\top + \boldsymbol{\Omega}]^{-1} \nu_t \quad (31)$$

where

$$\nu_t = \mathbf{y}_t - \boldsymbol{\Gamma}_t \boldsymbol{\Phi}_{t-1} \hat{\mathbf{x}}_{t-1|t-1} \quad (32)$$

We then have the following theorem.

Theorem 2. Suppose $\mathcal{Y}_t = \mathbf{Y}_t$. Denote by \mathcal{X}_t the set

$$\mathcal{X}_t = \{ \xi \in \mathbb{R}^n : (\xi - \hat{\mathbf{x}}_{t|t})^\top \mathbf{M}_{t|t}^{-1} (\xi - \hat{\mathbf{x}}_{t|t}) \leq \delta - \hat{\kappa}_t \} \quad (33)$$

and consider $t \in [0, T]$. Then the following statements hold

- 1) If $\hat{\kappa}_t \leq \delta$, $\forall t \in [0, T]$, then the expectation $\mathbb{E}_Q[\mathbf{X}_t|\mathcal{Y}_t]$ of (15) and (16) with $\mathbf{Q} \in \mathcal{Q}_{sec}$ belongs to the set \mathcal{X}_t .
- 2) The centroid of the set \mathcal{X}_t is the expected value $\mathbb{E}_P[\mathbf{X}_t|\mathcal{Y}_t]$ of the system (1) and (2) under P.
- 3) The variance of $\mathbb{E}_Q[\mathbf{X}_t|\mathcal{Y}_t] \in \mathcal{X}_t$ is $\mathbf{M}_{t|t}$.

The preceding theorem outlines a solution to the recursive estimation problem where the uncertainty in the system model is characterized by a change-of measure from P to Q where $\mathbf{Q} \in \mathcal{Q}_{sec}$ and \mathcal{Q}_{sec} is defined in Definition 1.

We also have a solution to the stochastic model validation problem since one can define the conditional probability that $\hat{\kappa}_t \leq \delta$ in terms of the underlying model and the constraint $\mathbf{Q} \in \mathcal{Q}_{sec}$. In other words, if $\hat{\kappa}_t > \delta$, for some $t \in [0, T]$ then there is some probability one could assign to the condition $\mathbf{Q} \notin \mathcal{Q}_{sec}$ or alternatively (but equivalently) there is some probability that the underlying model is not as described.

IV. ATTACK MODELS AND A DETECTION FRAMEWORK

Consider again the probability space $(\mathcal{S}, \mathcal{F}, \mathbb{P})$. Now consider the discrete-time stochastic system of the form

$$\mathbf{X}_{t+1} = \boldsymbol{\Phi}_t \mathbf{X}_t + \mathbf{E}_{t+1} \quad (34)$$

$$\mathbf{Z}_t = \boldsymbol{\Gamma}_t \mathbf{X}_t + \mathbf{N}_t + \alpha_t \quad (35)$$

where $\mathbf{X} \in \mathbb{R}^n$ denotes the state of the system and $\mathbf{Z} \in \mathbb{R}^l$ is the measured output of the system. Again, the notation $\mathbf{X} \in \mathbb{R}^n$ etc implies $\mathbf{X} : \mathcal{S} \rightarrow \mathbb{R}^n$. The random noise inputs to the system are $\mathbf{E} \in \mathbb{R}^n$ and $\mathbf{N} \in \mathbb{R}^l$.

The signal $\alpha_t \in \mathbb{R}^l$ is the so-called attack signal and is intended to distort the output of any estimator. In this case, α_t is a deterministic, but unknown, signal as expected. Consequently, the system (34) and (35) on the probability space $(\mathcal{S}, \mathcal{F}, \mathbb{Q})$ is simply

$$\mathbf{X}_{t+1} = \boldsymbol{\Phi}_t \mathbf{X}_t + \mathbf{U}_{t+1} + \zeta_{t+1} \quad (36)$$

$$\mathbf{Z}_t = \boldsymbol{\Gamma}_t \mathbf{X}_t + \mathbf{V}_t + \eta_t + \alpha_t \quad (37)$$

and again the measure we are working under should be clear from the context and notation.

Under P let \mathcal{Z}_t denote the completed σ -algebra generated by

$$\sigma(\{\mathbf{Z}_0, \dots, \mathbf{Z}_t\}) \quad (38)$$

and define a so-called innovation sequence $\{\tilde{\mathbf{Z}}_t\}$ by

$$\tilde{\mathbf{Z}}_t \triangleq \mathbf{Z}_t - \boldsymbol{\Gamma}_t \mathbb{E}_P[\mathbf{X}_t|\mathcal{Z}_{t-1}] \quad (39)$$

and let $\tilde{\mathcal{Z}}_t$ denote $\sigma(\{\tilde{\mathbf{Z}}_0, \dots, \tilde{\mathbf{Z}}_t\})$. Note that $\tilde{\mathbf{Z}}_t$ has mean α_t and covariance matrix

$$\boldsymbol{\Gamma}_t \text{Cov}_P(\mathbf{X}_t - \mathbb{E}_P[\mathbf{X}_t|\mathcal{Z}_{t-1}]) \boldsymbol{\Gamma}_t^\top + \boldsymbol{\Omega} \quad (40)$$

where $\text{Cov}_P(\mathbf{X}_t - \mathbb{E}_P[\mathbf{X}_t|\mathcal{Z}_{t-1}])$ denotes the covariance of $\mathbf{X}_t - \mathbb{E}_P[\mathbf{X}_t|\mathcal{Z}_{t-1}]$ under P.

Under Q define another so-called innovation sequence $\{\check{\mathbf{Z}}_t\}$ by

$$\check{\mathbf{Z}}_t \triangleq \mathbf{Z}_t - \boldsymbol{\Gamma}_t \mathbb{E}_Q[\mathbf{X}_t|\mathcal{Z}_{t-1}] \quad (41)$$

and let $\check{\mathcal{Z}}_t$ denote $\sigma(\{\check{\mathbf{Z}}_0, \dots, \check{\mathbf{Z}}_t\})$. Under Q it follows that $\check{\mathbf{Z}}_t$ has mean $(\mathbb{E}_Q[\eta_t] + \alpha_t)$ and covariance matrix

$$\boldsymbol{\Gamma}_t \text{Cov}_P(\mathbf{X}_t - \mathbb{E}_P[\mathbf{X}_t|\mathcal{Z}_{t-1}]) \boldsymbol{\Gamma}_t^\top + \boldsymbol{\Omega} \quad (42)$$

Now note that in general the innovations under P given by $\{\tilde{\mathbf{Z}}_t\}$ are not computable as the attack signal α_t is unknown. This is in contrast to the innovation sequence under P given by $\{\tilde{\mathbf{N}}_t\}$ which is typically computable given a realised sequence of measurements.

Consider now the test function

$$\Delta_P(k) = \begin{cases} 1 & \text{if } \mathbb{E}_P[\mathbf{Z}_t] \neq \mathbb{E}_P[\mathbf{Y}_t] \\ 0 & \text{if } \mathbb{E}_P[\mathbf{Z}_t] = \mathbb{E}_P[\mathbf{Y}_t] \end{cases} \quad (43)$$

and note the probability $\mathbb{P}(\Delta_P(k) \neq \Delta_Q(k)) = 0$ whenever $h(\mathbb{Q}||\mathbb{P}) \neq 0$.

Consider also the following test function

$$\Lambda_P(k) = \begin{cases} 1 & \text{if } \exists t \in \{0, \dots, k\} : \Delta_P(t) = 1 \\ 0 & \text{if } \Delta_P(t) = 0, \forall t \in \{0, \dots, k\} \end{cases} \quad (44)$$

such that if $\Lambda_P(k_0) = 1$ for some k_0 then $\Lambda_P(k) = 1$ for all $k \geq k_0$ while of course the same is not true for $\Delta_P(k)$. Again, we have the probability $P(\Lambda_P(k) \neq \Lambda_Q(k)) = 0$ whenever $h(Q||P) \neq 0$.

Of course, we do not know if our the output of our sensor in practice gives the following set of realized measurements $Y_t = \{y_0, \dots, y_t\}$ with $y_t \in \mathbb{R}^l$ or the following set of attack corrupted realized measurements $Z_t = \{z_0, \dots, z_t\}$ with $z_t \in \mathbb{R}^l$. It may, or is likely to, be combination of the two sets over a given time interval. As such, computing $\Delta_P(k)$ or $\Lambda_P(k)$ at every k is not possible. The focus of this work will be on estimating $\Lambda_P(k)$.

A. Attack Detection in Robust Stochastic Estimation

Recall that

$$\hat{\kappa}_t = \hat{\kappa}_{t-1} + \nu_t^\top [\Gamma_t \mathbf{M}_{t|t-1} \Gamma_t^\top + \Omega]^{-1} \nu_t \quad (45)$$

where

$$\nu_t = \mathbf{y}_t - \Gamma_t \Phi_{t-1} \mathbf{x}_{t-1|t-1} \quad (46)$$

and that $\hat{\kappa}_t < \delta$ is required to compute an ellipsoidal bound on the expected value $E_Q[\mathbf{X}_t|\mathcal{Y}_t]$ of (15) and (16) with $Q \in \mathcal{Q}_{sec}$. This bound is in the form of the set \mathcal{X}_t .

Theorem 3. Consider Definition 1 with some $\delta > 0$ and suppose the set of measurements is drawn from \mathcal{Y}_t ; i.e. there are no attacks on the sensor readings. Define

$$v_t = \begin{bmatrix} \mathbf{Y}_0 - \Gamma_0 \Phi_{-1} \mathbf{x}_{-1|t-1} \\ \vdots \\ \mathbf{Y}_t - \Gamma_t \Phi_{t-1} \mathbf{x}_{t-1|t-1} \end{bmatrix} \quad (47)$$

and

$$\Upsilon_t = \begin{bmatrix} \Gamma_0 \mathbf{M}_{0|t-1} \Gamma_0^\top + \Omega & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \Gamma_t \mathbf{M}_{t|t-1} \Gamma_t^\top + \Omega \end{bmatrix} \quad (48)$$

Then

$$P(\hat{\kappa}_t > \delta) = P(v_t^\top \Upsilon_t^{-1} v_t > \delta) \leq \frac{c}{\delta} \quad (49)$$

where $c = \text{trace}(\Upsilon_t^{-1} \Upsilon_t)$.

Proof of this theorem is a straightforward application of the Chebyshev inequality which puts a bound on the event probability that a random variable differs from its expected value by more than some specified amount.

In the situation where there is no attack on the system, δ is sufficiently large and $t < T$ then the probability that the expected value of the actual system $E_Q[\mathbf{X}_t|\mathcal{Y}_t]$ (considering (15) and (16)) for some $Q \in \mathcal{Q}_{sec}$ does not belong to the set \mathcal{X}_t is quite small. However, obviously the bound on $P(\hat{\kappa}_t > \delta)$ increases strictly with t . One would expect that as $t \rightarrow T$ (or more specifically, as the relative entropy introduced via the measure change approaches δ) then $P(\hat{\kappa}_t > \delta) \rightarrow 1$ at least loosely. Thus, for $t \geq T$ computing the set \mathcal{X}_t via δ should be increasingly difficult as the relative entropy between the nominal system and the uncertain system is likely increasing unmodelled.

Now we consider the effect of an attack sequence that distorts the measurement readings as previously specified. Thus, suppose one constructs \mathcal{X}_t as before but during the construction we substitute \mathbf{Y}_t with $\mathbf{Y}_t(1 - \Delta_P(t)) + \mathbf{Z}_t \Delta_P(t)$. Of course, the outcome of this substitution is unknown at the estimator which is derived based purely on the assumption that \mathbf{Y}_t is used. Then in practice one has

$$\tilde{\kappa}_t = \tilde{\kappa}_{t-1} + \hat{\nu}_t^\top [\Gamma_t \mathbf{M}_{t|t-1} \Gamma_t^\top + \Omega]^{-1} \hat{\nu}_t \quad (50)$$

where

$$\hat{\nu}_t = \mathbf{y}_t(1 - \Delta_P(t)) + \mathbf{z}_t \Delta_P(t) - \Gamma_t \Phi_{t-1} \mathbf{x}_{t-1|t-1} \quad (51)$$

but where it is still required that $\tilde{\kappa}_t < \delta$ in order to compute an ellipsoidal bound \mathcal{X}_t on the expected value $E_Q[\mathbf{X}_t|\mathcal{Y}_t]$ of (15) and (16) with $Q \in \mathcal{Q}_{sec}$.

Theorem 4. Consider Definition 1 with some $\delta > 0$ and suppose the set of measurements is drawn from $\mathbf{Y}_t(1 - \Delta_P(t)) + \mathbf{Z}_t \Delta_P(t)$; i.e. there may be attacks on the sensor readings. Let

$$\hat{\nu}_t = \begin{bmatrix} \mathbf{Y}_0(1 - \Delta_P(0)) + \mathbf{Z}_0 \Delta_P(0) - \Gamma_0 \Phi_{-1} \mathbf{x}_{-1|t-1} \\ \vdots \\ \mathbf{Y}_t(1 - \Delta_P(t)) + \mathbf{Z}_t \Delta_P(t) - \Gamma_t \Phi_{t-1} \mathbf{x}_{t-1|t-1} \end{bmatrix} \quad (52)$$

and define Υ_t as before. Then it follows that

$$P(\tilde{\kappa}_t > \delta) = P(\hat{\nu}_t^\top \Upsilon_t^{-1} \hat{\nu}_t > \delta) \leq \frac{c_1}{\delta - c_2} \quad (53)$$

with $\hat{\nu}_t$ now defining $\tilde{\kappa}_t$ and where $c_1 = \text{trace}(\Upsilon_t^{-1} \Upsilon_t)$ and

$$c_2 = \sum_{k=0}^t \alpha_k^\top [\Gamma_k \mathbf{M}_{k|k-1} \Gamma_k^\top + \Omega]^{-1} \alpha_k \quad (54)$$

Thus, it follows that the bound on $P(\tilde{\kappa}_t > \delta)$ increases strictly with each attack on the measurements. Define the attack detection estimator by

$$\hat{\Lambda}_P(k) = \begin{cases} 1 & \text{if } \tilde{\kappa}_t > \delta \\ 0 & \text{otherwise} \end{cases} \quad (55)$$

From the preceding two theorems it follows that the probability of detecting an attack (when there is an actual attack on the system) likely increases faster than the probability of generating a false alarm (whose bound grows only with time). Thus, while this attack detection scheme is likely to be conservative (especially as $t \ll T$) it is far less likely to generate a false alarm (particularly when $t \ll T$).

Note also that the bound on $P(\tilde{\kappa}_t > \delta)$ increases strictly with t (even when there is an attacker present) and thus the attacker will have an ever increasing challenge to attack the system and remain undetected.

V. CONCLUSION

The problem of recursive estimation and model validation for linear discrete-time systems with partial prior information was examined. An underlying linear discrete-time system is considered where the statistics of the driving noise is assumed to be known only partially; i.e. a class of noise inputs

is given from which the underlying actual noise is assumed to be chosen. A set-valued estimator is then derived and the conditional expectation is shown to belong to an ellipsoidal set consistent with the measurements and the underlying noise description. A method is also provided for estimating the consistency between the assumed model, knowledge on the partial prior noise statistics and the measured data. A group of attacking entities is then introduced with the goal of compromising the integrity of the state estimator by hijacking the sensor and distorting its output. It is shown that in order for the attack to go undetected, the distorted measurements need to be carefully designed and that given the model in this paper it is increasingly likely that an attack will be detected as time goes by.

VI. ACKNOWLEDGEMENTS

A.N. Bishop was supported by the NICTA, Control and Signal Processing Research Group and the NICTA, Canberra Research Laboratory. NICTA is funded by the Australian Government as represented by the Department of Broadband, Communications and the Digital Economy and the Australian Research Council through the ICT Centre of Excellence program. He is also supported by the US Air Force through AOARD (USAF-AOARD-10-4102) and the Australian Research Council (ARC) via a Discovery Early Career Researcher Award (DE-120102873).

REFERENCES

- [1] R.E. Kalman. A new approach to linear filtering and prediction problems. *Transaction of the ASME: Journal of Basic Engineering*, pages 35–45, March 1960.
- [2] T. Nishimura. On the a priori information in sequential estimation problems. *IEEE Transactions on Automatic Control*, 11(2):197–204, 1966.
- [3] R. Mehra. On the identification of variances and adaptive kalman filtering. *IEEE Transactions on Automatic Control*, 15(2):175–184, 1970.
- [4] S. Sangsuk-Iam and T.E. Bullock. Analysis of discrete-time Kalman filtering under incorrect noise covariances. *IEEE Transactions on Automatic Control*, 35(12):1304–1309, 1990.
- [5] J.L. Willems and F.M. Callier. Divergence of the stationary Kalman filter for correct and for incorrect noise variances. *IMA Journal of Mathematical Control and Information*, 9(1), 1992.
- [6] S. Sangsuk-Iam. Divergence of the discrete-time Kalman filter under incorrect noise covariances for linear periodic systems. In *Proceedings of the 1994 American Control Conference*, pages 1190–1194, 1994.
- [7] F. Schweppe. Recursive state estimation: unknown but bounded errors and system inputs. *IEEE Transactions on Automatic Control*, 13(1):22–28, 1968.
- [8] D. P. Bertsekas and I. B. Rhodes. Recursive state estimation for a set-membership description of uncertainty. *IEEE Transactions on Automatic Control*, 16(2):117–128, 1971.
- [9] J. Morris. The Kalman filter: A robust estimator for some classes of linear quadratic problems. *IEEE Transactions on Information Theory*, 22(5):526–534, 1976.
- [10] A.J. Krener. Kalman-bucy and minimax filtering. *IEEE Transactions on Automatic Control*, 25(2):291–292, 1980.
- [11] A.V. Savkin and I.R. Petersen. Recursive state estimation for uncertain systems with an integral quadratic constraint. *IEEE Transactions on Automatic Control*, 40(6):1080, 1995.
- [12] A.V. Savkin and I.R. Petersen. Model validation for robust control of uncertain systems with an integral quadratic constraint. *Automatica*, 32(4):603–606, 1996.
- [13] I.R. Petersen and A.V. Savkin. *Robust Kalman Filtering for Signals and Systems with Large Uncertainties*. Birkhauser, Boston, 1999.
- [14] I.R. Petersen, V.A. Ugrinovskii, and A.V. Savkin. *Robust Control Design Using H^∞ Methods*. Springer-Verlag, London, 2000.
- [15] A.V. Savkin and R.J. Evans. *Hybrid Dynamical Systems. Controller and Sensor Switching Problems*. Birkhauser, Boston, 2002.
- [16] U.D. Hanebeck, J. Horn, and G. Schmidt. On combining statistical and set-theoretic estimation. *Automatica*, 35(6):1101–1109, 1999.
- [17] F.K. Fletcher, M.S. Arulampalam, R.J. Evans, and W. Moran. Ellipsoidal set based tracking with nonlinear measurements. *IEE Proceedings on Radar, Sonar, and Navigation*, 152(5):335–344, October 2005.
- [18] B. Noack, V. Klumpp, D. Brunn, and U.D. Hanebeck. Nonlinear bayesian estimation with convex sets of probability densities. In *Proceedings of the 11th International Conference on Information Fusion*, 2008.
- [19] B. Noack, V. Klumpp, and U.D. Hanebeck. State estimation with sets of densities considering stochastic and systematic errors. In *Proceedings of the 12th International Conference on Information Fusion*, pages 1751–1758, 2009.
- [20] A. Dempster. Upper and lower probabilities induced by a multivalued mapping. *Classic Works of the Dempster-Shafer Theory of Belief Functions*, pages 57–72, 2008.
- [21] G. Dán and H. Sandberg. Stealth attacks and protection schemes for state estimators in power systems. In *Proceedings of the 1st IEEE International Conference on Smart Grid Communications*, pages 214–219, Gaithersburg, MD, USA, October 2010.
- [22] H. Sandberg, A. Teixeira, and K.H. Johansson. On security indices for state estimators in power networks. In *Proceedings of the 1st Workshop on Secure Control Systems*, Stockholm, Sweden, April 2010.
- [23] Y. Mo, E. Garone, A. Casavola, and B. Sinopoli. False data injection attacks against state estimation in wireless sensor networks. In *Proceedings of the 49th IEEE Conference on Decision and Control*, pages 5967–5972, Atlanta, GA, USA, December 2010.
- [24] S. Zheng, T. Jiang, and J.S. Baras. Robust state estimation under false data injection in distributed sensor networks. In *Proceedings of the 2010 IEEE Global Telecommunications Conference*, Miami, FL, USA, December 2010.
- [25] T. Kailath. The innovations approach to detection and estimation theory. *Proceedings of the IEEE*, 58(5):680–695, 1970.
- [26] S.K. Mitter. Filtering and stochastic control: a historical perspective. *IEEE Control Systems Magazine*, 16(3):67–76, 1996.
- [27] A.N. Bishop, P.N. Pathirana, and A.V. Savkin. Radar target tracking via robust linear filtering. *IEEE Signal Processing Letters*, 14(12):1028–1031, December 2007.
- [28] A.N. Bishop, P.N. Pathirana, and A.V. Savkin. Decentralized and robust target tracking with sensor networks. In *Proceedings of the 2008 IFAC World Congress*, Seoul, Korea, July 2008.
- [29] B.D.O. Anderson and J.B. Moore. *Optimal Filtering*. Prentice Hall, Englewood Cliffs, N.J., 1979.