

# Distributionally Robust State Estimation for Jump Linear Systems

Shixiong Wang

**Abstract**—In practice, the designed nominal model set for a jump (Markov) linear system might be uncertain: 1) Every candidate model might be inexact due to, e.g., mismatched modeling assumptions or model's identification errors; 2) The nominal model set might be incomplete (e.g., the true system has three operating modes but the designed model set includes only two of them). Moreover, the designed model transition probability matrix (TPM) might be uncertain: There is a discrepancy between the designed and true TPMs. Neglecting such model uncertainties and employing nominally optimal multi-model state estimators may cause significant performance losses. Therefore, this paper proposes a robust state estimation framework for jump linear systems that is insensitive to these three types of model uncertainty, technically by leveraging the distributionally robust optimization theory. Specifically, the model uncertainties are quantified by a collection of mixture distributions lying in a distributional ball centered at the nominal mixture distribution, where the nominal mixture distribution represents the nominal state-measurement distribution defined by nominal candidate models and their nominal model weights. Then, the robust state estimation is taken over the least-favorable distribution in the employed distributional ball. We show that the distributionally robust state estimation problem for jump linear systems can be reformulated into a tractable optimization equivalent such as a quadratic program or a positive semi-definite program, and in a special case, it can be analytically solved. Simulated and real-world experiments suggest that the proposed method is particularly useful when large model uncertainties exist.

**Index Terms**—Jump Linear Systems, Robust Filter, Distributional Robustness, Quadratic Program, Semi-Definite Program.

## I. INTRODUCTION

### A. Background

An actual physical plant or an information system may work in several different modes, and the modes may jump from one to another as time proceeds. For example, in multi-model target tracking [1], at each time step, a target can move according to any one of the following models [2]: the constant velocity (CV) model, the constant turn (CT) model, the constant acceleration (CA) model, the Singer model, etc. To track the evolution of this system, we are concerned with estimating its hidden state given measurements in the past. This paper considers Markov jump linear system models. The exactly optimal state estimation method for jump linear systems is computationally intractable because the number of the required state estimators grows exponentially as the time proceeds, and therefore, it is not implementable in practice [3],

[4]. Due to the intractability of the optimal state estimator, till now, the most popular method to handle the state estimation problem of jump linear systems is the interactive multiple model (IMM) filter, which **approximates** the filtered (i.e., posterior) state distribution using a limited number of Gaussian components [3], [5]. The IMM filter is pragmatically attractive because its computational burden is as low as the generalized pseudo-Bayesian estimator of first order (GPB1) but it has high performance as the generalized pseudo-Bayesian estimator of second order (GPB2) [3]. It is believed that the IMM filter can provide an excellent compromise between the filtering performance and the computational complexity [3], [6]–[10].

### B. Problem Statement

In practice, the IMM filter (and also other multi-model filters) faces the following limitations.

U1) The nominal model set might be uncertain; i.e., at some times, none of the nominal models in the nominal model set can exactly describe the true system dynamics. This can be understood from two aspects.

a) For every nominal model in the model set, it is an approximation to the true operating dynamics for the mode, and therefore, model mismatch exists. In multi-model target tracking, for instance, the nominal constant-turn (CT) model may be different from the true CT model because filter designers never exactly know the true turning rate, and as a result, a misidentified value of the turning rate may be used.

b) The nominal model set is not complete. For example, the number of the actual operating modes is larger than the size of the nominal model set. In multi-model target tracking, the target may move according to a great number of models [2], but filter designers may only use some of them (e.g., only CV, CA, and CT).

U2) The model transition probability matrix might be uncertain [7]–[9], [11], [12]. For example, in multi-model target tracking, model transition probabilities from one motion model to another may not be exactly known.

Hence, a state estimation framework for jump linear systems that is able to handle the listed two types of uncertainties in Items U1 and U2 is expected.

### C. Literature Review And Research Aims

The treatment frameworks for uncertainties in a single candidate model have been comprehensively surveyed and discussed in [13], [14], including the uncertainties in system

S. Wang is with the Institute of Data Science, National University of Singapore, Singapore 117602 (E-mail: s.wang@u.nus.edu).

This research is supported by the National Research Foundation Singapore and DSO National Laboratories under the AI Singapore Programme (AISG Award No: AISG2-RP-2020-018).

matrices and the uncertainties in noise distributions. When model sets are complete and TPMs are exact, these methods can be applied (or extended) for jump linear systems; cf., e.g., [15] and [16]. Therefore, this paper focuses on the review of treatment methods for the ad-hoc modeling uncertainties in jump linear systems, i.e., the incompleteness of nominal model sets and the uncertainties in TPMs. Since the incompleteness of nominal model sets has never been discussed in the literature, we pay attention to the uncertainties in TPMs. Several researches, e.g., [7]–[9], [11], [12], [17]–[19], have discussed the state estimation problem for jump linear systems when the model transition probability matrix is unknown or uncertain. These works can be categorized into two streams. The first stream aims at obtaining the accurate estimate of the **unknown** TPM using the Frequentist method [8], [12], [17] or the Bayesian method [7], [11], while the second stream tries to design robust state estimators that are insensitive to the **uncertain** TPM, e.g., the compensation-based method in [9] and the  $\mathcal{H}_\infty$  method in [18]–[20]. However,

- L1) The Frequentist method [8], [12], [17] and the Bayesian method [7], [11] assume that the true TPM is a time-invariant matrix. When the true TPM is (significantly) time-varying, these methods cannot provide the exact estimate of the true TPM anymore. In addition, even when the true TPM is time-invariant, a sufficiently long time horizon (i.e., sufficient measurements) is expected to estimate the TPM to a satisfactory level. Also, choosing satisfactory prior distributions for the unknown TPM and determining the parameters of the prior distributions are problem-dependent, and therefore, frustrating. Last but not least, albeit the true TPM is time-invariant, the estimate of the unknown TPM can hardly be exactly the same as its true value so that there still exists uncertainty (i.e., parameter identification error) in the estimated TPM;
- L2) The performance of the compensation-based method in [9] is problem-specific and it is desirable only when there exists an overwhelmingly dominating model at each time step (i.e., one of the model probabilities of candidate models is outstandingly large). The  $\mathcal{H}_\infty$  method in [18]–[20] requires that, for every candidate model, the noise sequences have finite energies, which implies that the power of the noises sequences vanish as time proceeds. Also, it requires the jump systems to be mean square stable. These two assumptions are obviously not always the case for the state estimation problem for general jump linear systems (e.g., in target tracking problems [2], [4], systems are usually non-stable).<sup>1</sup>

When the existence of model uncertainties is noticed but we have no specific information of where and how the uncertainties exist (e.g., whether system matrices, noise distributions, or TPMs are uncertain), the generic robust multi-model state estimators in [21], [22] are applicable. However, these generic methods tend to be overly conservative because they cannot

specifically respond to given types of uncertainty. To clarify further, when we exactly know that the model uncertainties exist only in TPMs, robust multi-model state estimators should not admit possible uncertainties in other components (e.g., system matrices and noise distributions) and only respond to uncertainties in TPMs.

Therefore, a unified multi-model state estimation framework for jump linear systems that is insensitive to the listed two types of uncertainty (i.e., Items U1 and U2) is expected. In particular, the new framework is supposed to flexibly respond to specified types of uncertainty.

#### D. Contributions

Following the conventional route [7]–[9], [11], [12], [17], this paper particularly robustifies the IMM filter. Specifically, a distributionally robust IMM filtering framework for jump linear systems that is insensitive to the two types of uncertainty listed in Subsection I-B is proposed; see (10), (11), and (12). Subsequently, the explicit optimization equivalents of the distributionally robust state estimation problem (11) subject to (12) are derived in Propositions 1 and 2. We show that these explicit optimization equivalents can be further reformulated into tractable optimization equivalents such as quadratic programs [see, e.g., (26)] or positive semi-definite programs [see, e.g., (34)], in different scenarios. Then, we show that these tractable reformulations can be efficiently solved using either the off-the-shelf solvers or the specifically-designed efficient algorithms (see, e.g., Proposition 4). Particularly, in a special case, the reformulated problem can be analytically solved; see Theorem 2. Finally, the distributionally robust IMM filter is summarized in Algorithm 1. Experiments suggest that the proposed method is particularly useful when complex and large model uncertainties exist in the nominal model; to be specific, for example, the method is especially suitable for tracking highly-maneuvering targets.

#### E. Notations

The space of all  $n$ -dimensional vectors is denoted by  $\mathbb{R}^n$ . Let  $\mathbb{P}_{\mathbf{x}}$  denote the distribution of the random vector  $\mathbf{x} \in \mathbb{R}^n$  (column by default). Let  $p_{\mathbf{x}}(\mathbf{x})$  denote the probability density (resp. mass) function of  $\mathbf{x}$  if  $\mathbf{x}$  is continuous (resp. discrete). Whenever it is clear from contexts, we use  $p(\mathbf{x})$  as a shorthand for  $p_{\mathbf{x}}(\mathbf{x})$ . The conditional distribution of  $\mathbf{x}$  given  $\mathbf{y} \in \mathbb{R}^m$  is denoted as  $\mathbb{P}_{\mathbf{x}|\mathbf{y}}$ . We use  $p_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y})$  to denote the conditional probability density (or mass) function of  $\mathbf{x}$  given  $\mathbf{y} = \mathbf{y}$ , shorted as  $p(\mathbf{x}|\mathbf{y})$ . Let  $\mathbb{E}\mathbf{x}$  denote the expectation of  $\mathbf{x}$  and  $\mathbb{E}(\mathbf{x}|\mathbf{y})$  the conditional expectation of  $\mathbf{x}$  given  $\mathbf{y}$ . The  $d$ -dimensional Gaussian distribution, parameterized by mean  $\mathbf{c}$  and covariance  $\Sigma$ , is denoted by  $\mathcal{N}_d(\mathbf{c}, \Sigma)$  and the corresponding Gaussian density function is denoted by  $\mathcal{N}_d(\mathbf{x}; \mathbf{c}, \Sigma)$ . Given an integer  $N$ , the running index set is defined as  $[N] := \{1, 2, \dots, N\}$ . Let  $\mathcal{Y}_k$  denote the measurement sequence up to and including the time  $k$ , i.e.,  $\mathcal{Y}_k := (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_k)$ . A realization of  $\mathcal{Y}_k$  is denoted as  $\mathbf{Y}_k := (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_k)$ . Let  $\mathbf{I}$  and  $\mathbf{0}$  denote the identity and the zero matrices with appropriate dimensions, respectively. We use  $\mathbf{M}^\top$  to denote the transpose of the matrix  $\mathbf{M}$ , and  $\text{Tr}[\mathbf{M}]$  its trace when

<sup>1</sup>However, the finite-energy assumption is reasonable for some state-feedback automatic control problems because some external disturbances are usually impulses or short-term perturbations. Besides, the mean-square-stability assumption is also reasonable for some state-feedback automatic control problems because systems can be stabilized by state feedback.

$M$  is square. Let  $\mathbb{S}^d$  denote the set of all  $d$ -dimensional symmetric matrices in  $\mathbb{R}^{d \times d}$ , and  $\mathbb{S}_+^d$  (resp.  $\mathbb{S}_{++}^d$ ) of all  $d$ -dimensional symmetric positive semi-definite (resp. positive definite) matrices in  $\mathbb{S}^d$ . If  $A, B \in \mathbb{S}^d$ ,  $A \succeq B$  (resp.  $A \succ B$ ) indicates that  $A - B \in \mathbb{S}_+^d$  (resp.  $A - B \in \mathbb{S}_{++}^d$ ). If  $S \in \mathbb{S}_+^d$ , let  $S^{1/2}$  be a square root of  $S$  (i.e.,  $S^{1/2}S^{1/2} = S$ ).

## II. PRELIMINARIES

### A. Bayesian Estimation Subject to Multiple Models

Suppose the random vectors  $\mathbf{x}$  and  $\mathbf{y}$  have finite second moments, and the joint distribution of them is  $\mathbb{P}_{\mathbf{x}, \mathbf{y}}$ . We are concerned with estimating the unobservable vector  $\mathbf{x}$  based on the model  $\mathbb{P}_{\mathbf{x}, \mathbf{y}}$  and the observation  $\mathbf{y}$ , in the minimum mean square error sense. In other words, we aim to find an estimator  $\hat{\mathbf{x}}$  of  $\mathbf{x}$ , which is a function  $\phi$  of  $\mathbf{y}$  [i.e.,  $\hat{\mathbf{x}} = \phi(\mathbf{y})$ ], such that

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\phi \in \mathcal{H}_{\mathbf{y}}} \operatorname{Tr} \mathbb{E}[\mathbf{x} - \phi(\mathbf{y})][\mathbf{x} - \phi(\mathbf{y})]^\top, \quad (1)$$

where the expectation is taken over the joint distribution  $\mathbb{P}_{\mathbf{x}, \mathbf{y}}$  and  $\mathcal{H}_{\mathbf{y}}$  contains all possible estimators of  $\mathbf{x}$  based on  $\mathbf{y}$ ; for more information, see [23]. As is well-known,  $\hat{\mathbf{x}} = \mathbb{E}(\mathbf{x}|\mathbf{y})$ . Sometimes, we are not sure of the exact form of  $\mathbb{P}_{\mathbf{x}, \mathbf{y}}$ . But we are confident that with probability  $\omega_j$ ,  $\mathbb{P}_{\mathbf{x}, \mathbf{y}}$  is of the form  $\mathbb{P}_{j, \mathbf{x}, \mathbf{y}}$  where  $\omega_j$ 's are weights and  $\sum_{j=1}^N \omega_j = 1$ . In other words,  $\mathbb{P}_{\mathbf{x}, \mathbf{y}}$  is a mixture of a set of distributions  $\{\mathbb{P}_{j, \mathbf{x}, \mathbf{y}}\}_{j=1,2,\dots,N}$  with mixing probabilities  $\{\omega_j\}_{j=1,2,\dots,N}$ , i.e.,  $\mathbb{P}_{\mathbf{x}, \mathbf{y}} = \sum_{j=1}^N \omega_j \mathbb{P}_{j, \mathbf{x}, \mathbf{y}}$ . As a result, given a measurement  $\mathbf{y}$ , the optimal estimate of  $\mathbf{x}$  is

$$\hat{\mathbf{x}} = \sum_{j=1}^N \mu_j \hat{\mathbf{x}}_j, \quad (2)$$

where  $\hat{\mathbf{x}}_j := \mathbb{E}(\mathbf{x}|\mathbf{y}, j) = \int \mathbf{x} p_j(\mathbf{x}|\mathbf{y}) d\mathbf{x}$ ,

$$\mu_j := p(j|\mathbf{y}) = \frac{\omega_j p_j(\mathbf{y})}{p(\mathbf{y})} = \frac{\omega_j p_j(\mathbf{y})}{\sum_{j=1}^N \omega_j p_j(\mathbf{y})}, \quad (3)$$

$p_j(\mathbf{x}|\mathbf{y}) := p(\mathbf{x}|\mathbf{y}, j)$  denotes the posterior distribution of  $\mathbf{x}$  given  $\mathbf{y}$  under the  $j^{\text{th}}$  model, and  $p_j(\mathbf{y}) := p(\mathbf{y}|j)$  denotes the likelihood of the  $j^{\text{th}}$  model under the measurement  $\mathbf{y}$ . Hence, by Bayes's rule,  $\mu_j \propto \omega_j p_j(\mathbf{y})$ . In other words,  $\omega_j$  can be understood as the prior model probability of the  $j^{\text{th}}$  model before observing  $\mathbf{y}$ , while  $\mu_j$  can be seen as the posterior model probability of the  $j^{\text{th}}$  model after observing  $\mathbf{y}$ . The corresponding posterior error covariance  $\mathbb{E}[(\mathbf{x} - \hat{\mathbf{x}})(\mathbf{x} - \hat{\mathbf{x}})^\top | \mathbf{y}]$  conditioned on  $\mathbf{y}$  equals to

$$\sum_{j=1}^N \mu_j [\mathbf{P}_j + (\hat{\mathbf{x}}_j - \hat{\mathbf{x}})(\hat{\mathbf{x}}_j - \hat{\mathbf{x}})^\top], \quad (4)$$

where  $\mathbf{P}_j := \mathbb{E}[(\mathbf{x} - \hat{\mathbf{x}}_j)(\mathbf{x} - \hat{\mathbf{x}}_j)^\top | j, \mathbf{y}]$  is the posterior error covariance of the  $j^{\text{th}}$  model given  $\mathbf{y}$ .

### B. Optimal Bayesian Estimation

The following fact is well established in applied statistics.

*Fact 1 ([23]):* The posterior mean  $\hat{\mathbf{x}} := \mathbb{E}(\mathbf{x}|\mathbf{y})$  solves both

$$\min_{\phi \in \mathcal{H}_{\mathbf{y}}} \operatorname{Tr} \mathbb{E}[\mathbf{x} - \phi(\mathbf{y})][\mathbf{x} - \phi(\mathbf{y})]^\top, \quad (5)$$

$$\text{and } \min_{\mathbf{a} \in \mathbb{R}^n} \operatorname{Tr} \mathbb{E}\{[\mathbf{x} - \mathbf{a}][\mathbf{x} - \mathbf{a}]^\top | \mathbf{y}\}, \quad \forall \mathbf{y} = \mathbf{y}, \quad (6)$$

where  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{y} \in \mathbb{R}^m$  are two random vectors;  $\mathbf{y}$  is a possible realization of  $\mathbf{y}$ . Note that given the realization  $\mathbf{y} = \mathbf{y}$ , the real number  $\mathbf{a} := \phi(\mathbf{y})$  is an *estimate* of  $\mathbf{x}$ .  $\square$

When  $\mathbf{y} = \mathbf{y}$  is specified, the optimal estimate of  $\mathbf{x}$  derived from (5) is  $\hat{\mathbf{x}} = \hat{\mathbf{x}}(\mathbf{y} = \mathbf{y}) = \mathbb{E}(\mathbf{x}|\mathbf{y} = \mathbf{y})$ ; i.e., the optimal estimate  $\hat{\mathbf{x}}$  is specified by the optimal estimator  $\hat{\mathbf{x}}$  with  $\mathbf{y}$  being replaced with  $\mathbf{y}$ . However, the optimal estimate derived from (6) given  $\mathbf{y} = \mathbf{y}$  is directly  $\hat{\mathbf{x}} = \mathbb{E}(\mathbf{x}|\mathbf{y})$ ; i.e., there is no optimal estimator directly associated with (6). We are interested in (6) because it is not always convenient to solve (5). Hence, when we obtain a measurement  $\mathbf{y}$ , we may directly solve (6). It is for this reason that in Subsection II-A, we work on a specified  $\mathbf{y}$  instead of the random vector  $\mathbf{y}$ . One may verify that, in Subsection II-A, deriving the closed-form expression of  $\mathbb{E}[(\mathbf{x} - \hat{\mathbf{x}})(\mathbf{x} - \hat{\mathbf{x}})^\top]$ , i.e.,  $\mathbb{E}_{\mathbf{y}} \mathbb{E}_{\mathbf{x}|\mathbf{y}}[(\mathbf{x} - \hat{\mathbf{x}})(\mathbf{x} - \hat{\mathbf{x}})^\top | \mathbf{y}]$  is extremely difficult. However, finding the closed-form expression of  $\mathbb{E}[(\mathbf{x} - \hat{\mathbf{x}})(\mathbf{x} - \hat{\mathbf{x}})^\top | \mathbf{y}]$  is relatively easy.

### C. Distributional Balls

A distributional ball centered at the reference distribution  $\bar{\mathbb{P}}_{\mathbf{x}}$  with radius  $\theta \geq 0$  is defined as  $\mathcal{F}_{\mathbf{x}}(\theta) := \{\mathbb{P}_{\mathbf{x}} | \Delta(\mathbb{P}_{\mathbf{x}}, \bar{\mathbb{P}}_{\mathbf{x}}) \leq \theta\}$ , where  $\Delta(\mathbb{P}_{\mathbf{x}}, \bar{\mathbb{P}}_{\mathbf{x}})$  is a proper statistical similarity measure (e.g., Wasserstein distance, Kullback–Leibler divergence, moment-based methods) between two distributions  $\mathbb{P}_{\mathbf{x}}$  and  $\bar{\mathbb{P}}_{\mathbf{x}}$  [13], [14];  $\theta$  is the size parameter of  $\mathcal{F}_{\mathbf{x}}(\theta)$ . Namely, a distributional ball is a collection of probability distributions that are close to  $\bar{\mathbb{P}}_{\mathbf{x}}$ . If  $\Delta$  is specified by a moment-based method [13, Eq. (21)], then we need more than one parameter to define the size of a distributional ball; in this case, the distributional ball is denoted as  $\mathcal{F}_{\mathbf{x}}(\boldsymbol{\theta})$  where a vector  $\boldsymbol{\theta}$  (called size parameter vector) is involved. Without loss of generality, this paper uses  $\mathcal{F}_{\mathbf{x}}(\boldsymbol{\theta})$  to denote a generic distributional ball.

## III. PROBLEM FORMULATION

We aim to estimate the unknown state  $\mathbf{x}_k$  of a jump linear Markov system

$$\begin{cases} \mathbf{x}_k = \mathbf{F}_{j,k-1} \mathbf{x}_{k-1} + \mathbf{G}_{j,k-1} \mathbf{w}_{j,k-1}, & j = 1, 2, \dots, N, \\ \mathbf{y}_k = \mathbf{H}_{j,k} \mathbf{x}_k + \mathbf{v}_{j,k}, \end{cases} \quad (7)$$

where  $k$  denotes the discrete time index;  $N$  is the size of the nominal model set;  $\mathbf{x}_k \in \mathbb{R}^n$  is the state vector;  $\mathbf{y}_k \in \mathbb{R}^m$  is the measurement vector;  $\mathbf{w}_{j,k-1} \in \mathbb{R}^p$ ,  $\mathbf{v}_{j,k} \in \mathbb{R}^m$  are the process noise and measurement noise of the  $j^{\text{th}}$  model, respectively. Typically, for every nominal linear system in the model set (i.e., for every  $j = 1, 2, \dots, N$ ), the following properties are assumed to be satisfied [24]–[26]: 1) for all  $k$ ,  $\mathbf{x}_k$ ,  $\mathbf{y}_k$ ,  $\mathbf{w}_{j,k}$ , and  $\mathbf{v}_{j,k}$  have finite second moments; 2)  $\mathbf{x}_0 \sim \mathcal{N}_n(\bar{\mathbf{x}}_0, \mathbf{M}_0)$ , and for all  $k$ ,  $\mathbf{w}_{j,k} \sim \mathcal{N}_p(\boldsymbol{\mu}_{j,k}^w, \mathbf{Q}_{j,k})$  and  $\mathbf{v}_{j,k} \sim \mathcal{N}_m(\boldsymbol{\mu}_{j,k}^v, \mathbf{R}_{j,k})$ ; 3) for any  $k_1 \neq k_2$ ,  $\mathbf{w}_{j,k_1}$  and  $\mathbf{x}_0$  are uncorrelated, so are  $\mathbf{v}_{j,k_1}$  and  $\mathbf{x}_0$ ,  $\mathbf{w}_{j,k_1}$  and  $\mathbf{w}_{j,k_2}$ , and  $\mathbf{v}_{j,k_1}$  and  $\mathbf{v}_{j,k_2}$ ; for any  $k_1, k_2$ ,  $\mathbf{v}_{j,k_1}$  and  $\mathbf{w}_{j,k_2}$  are uncorrelated; 4) the involved parameters  $\bar{\mathbf{x}}_0$ ,  $\mathbf{M}_0$ ,  $\boldsymbol{\mu}_{j,k}^w$ ,  $\boldsymbol{\mu}_{j,k}^v$ ,  $\mathbf{Q}_{j,k}$ ,  $\mathbf{R}_{j,k}$ ,  $\mathbf{F}_{j,k-1}$ ,  $\mathbf{G}_{j,k-1}$ , and  $\mathbf{H}_{j,k}$  are exactly known, and typically  $\boldsymbol{\mu}_{j,k}^w$  and  $\boldsymbol{\mu}_{j,k}^v$  are the zero vectors. Note that for a true operating jump system, at every time  $k$ , the true dominating working mode  $j$  is unknown.

The jump linear system (7) is called a hybrid linear system because we can treat the operating but unknown mode  $j_k$  at the time  $k$  as a discrete system state so that the augmented state vector  $\{j_k, \mathbf{x}_k\}$  is a hybrid state vector consisting of both a discrete state variable and a continuous state vector. Usually, the evolution of the discrete state  $j_k$  (i.e., the model transition process) is modeled by a  $N$ -state homogeneous Markov chain [3], and the model transition probability matrix (TPM) is  $\Pi := \{\pi_{ij}\}_{i,j=1,2,\dots,N}$  where  $\pi_{ij}$  denotes the probability that the system's operating mode jumps from the  $i^{\text{th}}$  model at the time  $k-1$  to the  $j^{\text{th}}$  model at the time  $k$ . From the viewpoint of Bayesian statistical signal processing, the state estimation problem for jump linear systems can be stated as finding the posterior (or filtered) state distribution of  $\{j_k, \mathbf{x}_k\}$ , i.e.,  $p(j_k, \mathbf{x}_k | \mathbf{Y}_k)$ , at the time  $k$ , based on the system model (7), the TPM  $\Pi$ , and the past measurements  $\mathbf{Y}_k := (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_k)$ .

The nominal system (7) defines two discrete-time stochastic processes  $\{\mathbf{x}_k\}$  and  $\{\mathbf{y}_k\}$ ,  $k = 1, 2, \dots$ . Suppose the nominal joint state-measurement distribution defined by the nominal system model (7) is  $\mathbb{P}_{\mathbf{x}_k, \mathbf{y}_k}$ . We would like to solve the following optimization problem

$$\min_{\phi \in \mathcal{H}_{\mathcal{Y}_k}} \text{Tr} \mathbb{E}[\mathbf{x}_k - \phi(\mathcal{Y}_k)][\mathbf{x}_k - \phi(\mathcal{Y}_k)]^\top, \quad (8)$$

where the expectation is taken over the joint distribution  $\mathbb{P}_{\mathbf{x}_k, \mathbf{y}_k}$ ;  $\mathcal{H}_{\mathcal{Y}_k}$  contains all possible estimators of  $\mathbf{x}_k$  based on the measurement set  $\mathcal{Y}_k$ ;  $\phi$  is called an estimator and the one solving (8) is the optimal estimator. One may easily verify by contradiction that the optimal solution to (8) is unique. The optimal estimator of  $\mathbf{x}_k$  in this minimum mean square error sense is  $\mathbb{E}(\mathbf{x}_k | \mathcal{Y}_k) \in \mathcal{H}_{\mathcal{Y}_k}$ . Since  $\mathbb{P}_{\mathbf{x}_k, \mathbf{y}_k}$  is not Gaussian due to the multi-mode property of (7),  $\mathbb{E}(\mathbf{x}_k | \mathcal{Y}_k)$  cannot be of a linear form because a Gaussian mixture is no longer Gaussian; cf. (2) and (3) where  $\hat{\mathbf{x}}$  is no longer linear in  $\mathbf{y}$ .

If the actual system dynamics deviates from the nominal model (7), the true joint state-measurement distribution  $\mathbb{P}_{\mathbf{x}_k, \mathbf{y}_k}$  will deviate from the nominal  $\mathbb{P}_{\mathbf{x}_k, \mathbf{y}_k}$ . As a result, a robust state estimation solution that is insensitive to the model deviations needs to be designed. The distributionally robust counterpart of (8) can be written as

$$\min_{\phi \in \mathcal{H}_{\mathcal{Y}_k}} \max_{\mathbb{P} \in \mathcal{F}_{\mathbf{x}_k, \mathbf{y}_k}(\boldsymbol{\theta})} \text{Tr} \mathbb{E}[\mathbf{x}_k - \phi(\mathcal{Y}_k)][\mathbf{x}_k - \phi(\mathcal{Y}_k)]^\top, \quad (9)$$

where the expectation is taken over  $\mathbb{P}_{\mathbf{x}_k, \mathbf{y}_k}$  and  $\mathcal{F}_{\mathbf{x}_k, \mathbf{y}_k}(\boldsymbol{\theta})$  is the associated **ambiguity set** consisting of all possible joint distributions  $\mathbb{P}_{\mathbf{x}_k, \mathbf{y}_k}$  that lie in a distributional ball centered at the nominal distribution  $\mathbb{P}_{\mathbf{x}_k, \mathbf{y}_k}$  with size parameter vector  $\boldsymbol{\theta}$ .

Since state estimation problems are real-time problems: the measurements  $\mathbf{y}_k$  arrives sequentially and the optimal estimator operates along the discrete time in a recursive way [27], we instead solve a time-incremental [28] problem

$$\min_{\phi \in \mathcal{H}_{\mathcal{Y}_k}} \max_{\mathbb{P} \in \mathcal{F}_{\mathbf{x}_k, \mathbf{y}_k | \mathcal{Y}_{k-1}}(\boldsymbol{\theta})} \text{Tr} \mathbb{E}\{[\mathbf{x}_k - \phi(\mathbf{y}_k)][\mathbf{x}_k - \phi(\mathbf{y}_k)]^\top | \mathcal{Y}_{k-1}\}, \quad (10)$$

where the expectation is taken over  $\mathbb{P}_{\mathbf{x}_k, \mathbf{y}_k | \mathcal{Y}_{k-1}}$  and the ambiguity set  $\mathcal{F}_{\mathbf{x}_k, \mathbf{y}_k | \mathcal{Y}_{k-1}}(\boldsymbol{\theta})$  is constructed around  $\mathbb{P}_{\mathbf{x}_k, \mathbf{y}_k | \mathcal{Y}_{k-1}}$ , i.e., the nominal conditional joint state-measurement distribution given the previous measurement sequence. Note that in (10), the space of  $\phi$  is only defined by  $\mathbf{y}_k$  instead of

$\mathcal{Y}_k$ . To solve (10), we need to first design proper forms for  $\mathcal{F}_{\mathbf{x}_k, \mathbf{y}_k | \mathcal{Y}_{k-1}}(\boldsymbol{\theta})$ , and then find the explicit optimization equivalent(s) of (10) so that it can be efficiently solved.

Therefore, at each time step  $k$ , we are inspired to **first** study a distributionally robust Bayesian estimation problem with multiple nominal models

$$\min_{\phi \in \mathcal{H}_{\mathcal{Y}}} \max_{\mathbb{P} \in \mathcal{F}_{\mathbf{x}, \mathbf{y}}(\boldsymbol{\theta})} \text{Tr} \mathbb{E}[\mathbf{x} - \phi(\mathbf{y})][\mathbf{x} - \phi(\mathbf{y})]^\top \quad (11)$$

subject to the multiple nominal joint state-measurement distributions  $\{\bar{\mathbb{P}}_{j, \mathbf{x}, \mathbf{y}}\}_{j \in [N]}$ , the nominal prior model probabilities  $\{\bar{\omega}_j\}_{j \in [N]}$ , and a properly constructed ambiguity set  $\mathcal{F}_{\mathbf{x}, \mathbf{y}}(\boldsymbol{\theta})$  that is characterized by  $\{\bar{\mathbb{P}}_{j, \mathbf{x}, \mathbf{y}}\}_{j \in [N]}$  and  $\{\bar{\omega}_j\}_{j \in [N]}$ . The subscript  $k$  (i.e., discrete time index) is dropped to avoid notational clutter. **Then**, by identifying the joint distribution of  $(\mathbf{x}_k, \mathbf{y}_k)$  conditioned on  $\mathcal{Y}_{k-1}$ , we can solve (10).

*Remark 1:* Model (11) is called the distributionally robust counterpart of (1). Note that robust counterparts of the batch and the recursive formulations are not equivalent: The former (9) is called robust filtering without commitment while the latter (10) is called robust filtering under commitment [28]. However, the robust counterpart for the recursive formulation outperforms that for the batch formulation because the latter tends to consume all the prescribed robustness budget at the first several time steps (i.e., as the time proceeds, the robust batch formulation would reduce to the usual non-robust version); see [29, Appendix B].  $\square$

#### IV. DISTRIBUTIONALLY ROBUST BAYESIAN ESTIMATION SUBJECT TO MULTIPLE NOMINAL MODELS

In this section, we first design a proper ambiguity set  $\mathcal{F}_{\mathbf{x}, \mathbf{y}}(\boldsymbol{\theta})$  that is compatible with the research aims raised in Subsection I-C, and then solve the distributionally robust Bayesian estimation problem (11). Because we desire the flexibility that responds to specified types of uncertainty, in order to take into account the individual model uncertainties in every nominal candidate model and/or the uncertainties in nominal model weights,<sup>2</sup> a suitable ambiguity set  $\mathcal{F}_{\mathbf{x}, \mathbf{y}}(\boldsymbol{\theta})$  can be constructed in (12), shown at the top of the next page, where  $\boldsymbol{\theta} := [\theta_0, \theta_1, \theta_2, \dots, \theta_N]^\top$ ,  $\mathcal{M}(\mathbb{R}^d)$  denotes all probability distributions on  $\mathbb{R}^d$ , and  $\mathcal{P}(\mathbb{R}^N)$  denotes all  $N$ -length discrete distributions;  $\bar{\mathbb{P}}_{j, \mathbf{x}, \mathbf{y}}$  represents the  $j^{\text{th}}$  nominal joint state-measurement distribution;  $\bar{\omega} := (\bar{\omega}_1, \bar{\omega}_2, \dots, \bar{\omega}_N)^\top$  and  $\bar{\omega}_j$  denotes the nominal prior model probability of the  $j^{\text{th}}$  model; for all  $j \in [N]$ ,  $\Delta_j(\cdot, \cdot)$  represents the statistical similarity measure between the possibly true joint state-measurement distribution  $\mathbb{P}_{j, \mathbf{x}, \mathbf{y}}$  and its nominal distribution  $\bar{\mathbb{P}}_{j, \mathbf{x}, \mathbf{y}}$ , and  $\Delta_0(\cdot, \cdot)$  the statistical similarity measure between the possibly true prior model probability  $\omega$  and its nominal value  $\bar{\omega}$ ;  $\theta_0$  and  $\theta_j$  are the radii of the distributional sets defined by  $\Delta_0(\cdot, \cdot)$  and  $\Delta_j(\cdot, \cdot)$ 's, respectively. The ambiguity set in (12) inherently allows the flexibility of responding to specified types of uncertainty. For example, when we set  $\theta_0$  to be non-zero and all other  $\theta$ s (from  $\theta_1$  to  $\theta_N$ ) to be zeros, we consider only the uncertainties in TPMs (and therefore in model weights), and all candidate models are believed to be

<sup>2</sup>N.B.: Uncertainties in TPMs lead to uncertainties in model weights.

$$\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta}) = \left\{ \mathbb{P}_{\mathbf{x},\mathbf{y}} = \sum_{j=1}^N \omega_j \mathbb{P}_{j,\mathbf{x},\mathbf{y}} \left| \begin{array}{l} \mathbb{P}_{j,\mathbf{x},\mathbf{y}} \in \mathcal{M}(\mathbb{R}^n \times \mathbb{R}^m), \quad \forall j \in [N] \\ \boldsymbol{\omega} \in \mathcal{P}(\mathbb{R}^N) \\ \Delta_j(\mathbb{P}_{j,\mathbf{x},\mathbf{y}}, \bar{\mathbb{P}}_{j,\mathbf{x},\mathbf{y}}) \leq \theta_j, \quad \forall j \in [N] \\ \Delta_0(\boldsymbol{\omega}, \bar{\boldsymbol{\omega}}) \leq \theta_0 \end{array} \right. \right\}, \quad (12)$$

exact. For another example, supposing  $N = 2$ , if we set  $\theta_1$  to be non-zero and  $\theta_0$  and  $\theta_2$  to be zeros, we consider only the model uncertainties in the first nominal candidate model, and the TPM and the second candidate model are believed to be exact. Note that every element in  $\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$  is a mixture distribution, including the least-favorable distribution that solves the inner maximization problem in (11). If  $N = 1$ , we have  $\omega_1 = \bar{\omega}_1 = 1$  and  $\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$  reduces to

$$\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta}) = \{ \mathbb{P}_{\mathbf{x},\mathbf{y}} \in \mathcal{M}(\mathbb{R}^n \times \mathbb{R}^m) \mid \Delta(\mathbb{P}_{\mathbf{x},\mathbf{y}}, \bar{\mathbb{P}}_{\mathbf{x},\mathbf{y}}) \leq \theta \},$$

which is a special case discussed in [13].

Next, we find the tractable reformulation(s) of the distributionally robust Bayesian estimation problem (11) subject to the multi-model ambiguity set (12).

Consider the max-min problem induced by (11):

$$\max_{\mathbb{P} \in \mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})} \min_{\phi \in \mathcal{H}_{\mathbf{y}}} \text{Tr} \mathbb{E}[\mathbf{x} - \phi(\mathbf{y})][\mathbf{x} - \phi(\mathbf{y})]^\top. \quad (13)$$

*Lemma 1:* The min-max problem (11) and the max-min problem (13) is equivalent. To be specific, the solution solving (13) also solves (11), and vice versa. To clarify further, letting  $V(\phi, \mathbb{P}) := \text{Tr} \mathbb{E}[\mathbf{x} - \phi(\mathbf{y})][\mathbf{x} - \phi(\mathbf{y})]^\top$ , and supposing  $\phi^*$  and  $\mathbb{P}^*$  solve the max-min problem (13), then  $(\phi^*, \mathbb{P}^*)$  forms a saddle point of the objective function  $V(\phi, \mathbb{P})$ , i.e.,

$$\min_{\phi \in \mathcal{H}_{\mathbf{y}}} V(\phi, \mathbb{P}^*) = V(\phi^*, \mathbb{P}^*) = \max_{\mathbb{P} \in \mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})} V(\phi^*, \mathbb{P}). \quad (14)$$

*Proof:* The Gaussianity assumption in [14, Thm. 7] and [30, Thm. 9] can be dropped without changing the statements therein because, for any type of joint distribution  $\mathbb{P}_{\mathbf{x},\mathbf{y}}$  (not limited to a Gaussian), the unique optimal state estimator in the minimum mean square error sense is the posterior mean  $\mathbb{E}(\mathbf{x}|\mathbf{y})$ . As a result, this lemma is immediate following the proofs of [14, Thm. 7] and [30, Thm. 9].  $\square$

The max-min problem (13) is easier to solve than the original min-max problem (11) because for every  $\mathbb{P} \in \mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$ , we can find the unique associated optimal estimator in  $\mathcal{H}_{\mathbf{y}}$ . For every  $j \in [N]$ , supposing the nominal joint state-noise distribution is  $\bar{\mathbb{P}}_{j,\mathbf{x},\mathbf{v}_j} := \mathcal{N}_{n+m} \left( \begin{bmatrix} \bar{\mathbf{x}}_j \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{M}_j & \mathbf{0} \\ \mathbf{0} & \mathbf{R}_j \end{bmatrix} \right)$ , we have the nominal joint state-measurement distribution

$$\bar{\mathbb{P}}_{j,\mathbf{x},\mathbf{y}} := \mathcal{N}_{n+m} \left( \begin{bmatrix} \bar{\mathbf{x}}_j \\ \mathbf{H}_j \bar{\mathbf{x}}_j \end{bmatrix}, \begin{bmatrix} \mathbf{M}_j & \mathbf{M}_j \mathbf{H}_j^\top \\ \mathbf{H}_j \mathbf{M}_j & \mathbf{H}_j \mathbf{M}_j \mathbf{H}_j^\top + \mathbf{R}_j \end{bmatrix} \right), \quad (15)$$

which is induced by the linear observation system  $\mathbf{y} = \mathbf{H}_j \mathbf{x} + \mathbf{v}_j$ . Accordingly, given the measurement  $\mathbf{y}$ , we can obtain the nominal state estimate  $\hat{\mathbf{x}}_j$  and the nominal estimation error covariance  $\bar{\mathbf{P}}_j$ ; see Appendix A. For every candidate distribution  $\mathbb{P}_{j,\mathbf{x},\mathbf{y}}$  in  $\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$ , we assume that its component distributions  $\mathbb{P}_{j,\mathbf{x},\mathbf{y}}, \forall j \in [N]$  are also Gaussian. If the possible joint state-noise distribution is  $\bar{\mathbb{P}}_{j,\mathbf{x},\mathbf{v}_j} :=$

$\mathcal{N}_{n+m} \left( \begin{bmatrix} \mathbf{c}_{j,\mathbf{x}} \\ \mathbf{c}_{j,\mathbf{v}} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{j,\mathbf{x}} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{j,\mathbf{v}} \end{bmatrix} \right)$ , we can obtain the joint state-measurement distribution

$$\mathbb{P}_{j,\mathbf{x},\mathbf{y}} := \mathcal{N}_{n+m}(\mathbf{c}_j, \boldsymbol{\Sigma}_j) \quad (16)$$

where  $\boldsymbol{\Sigma}_j := \begin{bmatrix} \boldsymbol{\Sigma}_{j,\mathbf{x}} & \boldsymbol{\Sigma}_{j,\mathbf{x}} \mathbf{H}_j^\top \\ \mathbf{H}_j \boldsymbol{\Sigma}_{j,\mathbf{x}} & \mathbf{H}_j \boldsymbol{\Sigma}_{j,\mathbf{x}} \mathbf{H}_j^\top + \boldsymbol{\Sigma}_{j,\mathbf{v}} \end{bmatrix}$  and  $\mathbf{c}_j := \begin{bmatrix} \mathbf{c}_{j,\mathbf{x}} \\ \mathbf{H}_j \mathbf{c}_{j,\mathbf{x}} + \mathbf{c}_{j,\mathbf{v}} \end{bmatrix}$ . Accordingly, given the measurement  $\mathbf{y}$ , the state estimate  $\hat{\mathbf{x}}_j$  and the estimation error covariance  $\mathbf{P}_j$  associated with  $\mathbb{P}_{j,\mathbf{x},\mathbf{y}}$  can be obtained; see Appendix A.

From (2) and (4), for a possible mixture distribution  $\mathbb{P}_{\mathbf{x},\mathbf{y}} \in \mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$ , if  $\mathbf{y} = \mathbf{y}$  is specified, we have the associated optimal posterior estimate

$$\hat{\mathbf{x}} = \sum_{j=1}^N \mu_j \cdot \hat{\mathbf{x}}_j, \quad (17)$$

and the corresponding estimation error covariance

$$\mathbf{P} = \sum_{j=1}^N \mu_j \cdot \{ \mathbf{P}_j + (\hat{\mathbf{x}}_j - \hat{\mathbf{x}})(\hat{\mathbf{x}}_j - \hat{\mathbf{x}})^\top \}, \quad (18)$$

where  $\hat{\mathbf{x}}_j$ ,  $\mathbf{P}_j$ , and  $\mu_j$  are the optimal posterior estimate of the state  $\mathbf{x}$ , the posterior error covariance, and the posterior model probability, corresponding to the  $j^{\text{th}}$  nominal model, respectively. As explained in Subsection II-B, to simplify the problem-solving procedure and without loss of optimality, we work directly on the case where  $\mathbf{y}$  has been specified; i.e., (6).

Consequently, the explicit optimization equivalent of the distributionally robust Bayesian estimation problem (11) subject to the ambiguity set (12) can be given below.

*Proposition 1:* Consider the distributionally robust Bayesian estimation problem (11). Suppose that

- 1) The ambiguity set  $\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$ , which is a collection of mixture distributions, is defined in (12);
- 2) Every candidate distribution inside  $\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$  is a  $N$ -component Gaussian mixture, and each Gaussian component is defined as (16), for every  $j \in [N]$ ;
- 3) The nominal distribution, which is the center of  $\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$ , is a  $N$ -component Gaussian mixture, and each Gaussian component is defined as (15), for every  $j \in [N]$ ;
- 4) A measurement  $\mathbf{y} = \mathbf{y}$  is specified.

Then (11) can be explicitly reformulated into

$$\max_{\boldsymbol{\Theta}} \text{Tr} \sum_{j=1}^N \mu_j \cdot \{ \mathbf{P}_j + (\hat{\mathbf{x}}_j - \hat{\mathbf{x}})(\hat{\mathbf{x}}_j - \hat{\mathbf{x}})^\top \}, \quad (19)$$

where  $\boldsymbol{\Theta} := \{ \omega_j, \mathbf{c}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{x}}, \mathbf{c}_{j,\mathbf{v}}, \boldsymbol{\Sigma}_{j,\mathbf{v}} \}_{j \in [N]}$ ; the feasible region of  $\boldsymbol{\Theta}$  is implicitly defined by  $\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$  because every  $\mathbb{P}_{\mathbf{x},\mathbf{y}}$  in  $\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$  is parameterized by  $\boldsymbol{\Theta}$ . Note that for every  $j \in [N]$ ,  $\mu_j$  is defined by  $\omega_j$  through (3), and  $\hat{\mathbf{x}}_j$  and  $\mathbf{P}_j$  are defined by  $\mathbf{c}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{x}}, \mathbf{c}_{j,\mathbf{v}}$ , and  $\boldsymbol{\Sigma}_{j,\mathbf{v}}$  (cf. Appendix A), respectively;  $\hat{\mathbf{x}}$  is defined in (17).

*Proof:* See Appendix B.  $\square$

Compared with the original distributionally robust Bayesian estimation problem (11), which is extremely abstract, the reformulated optimization equivalent (19) is explicit and specific.

Hence, when the ambiguity set  $\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$  is designed, (19) can be explicitly solved over  $\boldsymbol{\Theta}$ . However, (19) is still a difficult problem because the objective is highly nonlinear in  $\boldsymbol{\Theta}$ .

For the convenience of later analysis, we give a compact reformulation of (19).

*Proposition 2:* Problem (19) can be compactly rewritten as

$$\max_{\boldsymbol{\Theta}} -\boldsymbol{\mu}^\top \mathbf{A}\boldsymbol{\mu} + \mathbf{b}^\top \boldsymbol{\mu}, \quad (20)$$

where

$$\mathbf{A} := \begin{bmatrix} \hat{\mathbf{x}}_1^\top \hat{\mathbf{x}}_1 & \hat{\mathbf{x}}_1^\top \hat{\mathbf{x}}_2 & \hat{\mathbf{x}}_1^\top \hat{\mathbf{x}}_3 & \cdots & \hat{\mathbf{x}}_1^\top \hat{\mathbf{x}}_N \\ \hat{\mathbf{x}}_2^\top \hat{\mathbf{x}}_1 & \hat{\mathbf{x}}_2^\top \hat{\mathbf{x}}_2 & \hat{\mathbf{x}}_2^\top \hat{\mathbf{x}}_3 & \cdots & \hat{\mathbf{x}}_2^\top \hat{\mathbf{x}}_N \\ \hat{\mathbf{x}}_3^\top \hat{\mathbf{x}}_1 & \hat{\mathbf{x}}_3^\top \hat{\mathbf{x}}_2 & \hat{\mathbf{x}}_3^\top \hat{\mathbf{x}}_3 & \cdots & \hat{\mathbf{x}}_3^\top \hat{\mathbf{x}}_N \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \hat{\mathbf{x}}_N^\top \hat{\mathbf{x}}_1 & \hat{\mathbf{x}}_N^\top \hat{\mathbf{x}}_2 & \hat{\mathbf{x}}_N^\top \hat{\mathbf{x}}_3 & \cdots & \hat{\mathbf{x}}_N^\top \hat{\mathbf{x}}_N \end{bmatrix}, \quad (21)$$

and

$$\mathbf{b} := \begin{bmatrix} \text{Tr}[\mathbf{P}_1] + \hat{\mathbf{x}}_1^\top \hat{\mathbf{x}}_1 \\ \text{Tr}[\mathbf{P}_2] + \hat{\mathbf{x}}_2^\top \hat{\mathbf{x}}_2 \\ \vdots \\ \text{Tr}[\mathbf{P}_N] + \hat{\mathbf{x}}_N^\top \hat{\mathbf{x}}_N \end{bmatrix}. \quad (22)$$

Since  $\mathbf{A} \succeq \mathbf{0}$ , the objective of (20) is concave in  $\boldsymbol{\mu}$ .

*Proof:* The objective of (19) can be rearranged into  $-\text{Tr}[(\sum_{j=1}^N \mu_j \cdot \hat{\mathbf{x}}_j)(\sum_{j=1}^N \mu_j \cdot \hat{\mathbf{x}}_j)^\top] + \sum_{j=1}^N \mu_j \cdot \text{Tr}[\mathbf{P}_j + \hat{\mathbf{x}}_j \hat{\mathbf{x}}_j^\top]$ . By defining  $\mathbf{A}$ ,  $\mathbf{b}$ , we have (20).  $\square$

The reformulation from (19) to (20) is important for later analysis because the objective of the former is cubic in  $\boldsymbol{\mu}$  but that of the latter is quadratic in  $\boldsymbol{\mu}$ . Cubic programming is difficult because convexity (or concavity) is not naturally implied. However, generally speaking, (20) is still complicated to be solved because  $\boldsymbol{\mu}$ ,  $\mathbf{A}$ , and  $\mathbf{b}$  are all defined by the decision variables  $\boldsymbol{\Theta}$  in a highly nonlinear manner.

In the following, we explicitly solve (19) or (20), whichever is easier under investigated conditions.

*A. When  $\theta_0 \neq 0$  But  $\theta_j = 0, \forall j \in [N]$*

We first study a special case where  $\theta_j = 0, \forall j = 1, 2, \dots, N$ . Namely, we assume that the nominal model set is correct and able to exactly describe every true dynamics of the actual system. However, the prior model probabilities  $\omega_j$ 's are uncertain due, e.g., to the uncertain model transition probability matrix, to the uncertain initial model probability at the time  $k = 0$ , or to the method uncertainty (of, e.g., the IMM filter) in approximating posterior state distributions.

When  $\theta_j = 0, \forall j \in [N]$ ,  $\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$  in (12) reduces to

$$\mathcal{F}_{\mathbf{x},\mathbf{y}}(\theta_0) = \left\{ \mathbb{P}_{\mathbf{x},\mathbf{y}} = \sum_{j=1}^N \omega_j \bar{\mathbb{P}}_{j,\mathbf{x},\mathbf{y}} \mid \begin{array}{l} \boldsymbol{\omega} \in \mathcal{P}(\mathbb{R}^N) \\ \Delta_0(\boldsymbol{\omega}, \bar{\boldsymbol{\omega}}) \leq \theta_0 \end{array} \right\}. \quad (23)$$

It means that the nominal model weight vector  $\boldsymbol{\omega}$  is uncertain, but the nominal model set is exact. Since only  $\boldsymbol{\mu}$  is related to  $\boldsymbol{\omega}$ , and  $\mathbf{A}$  and  $\mathbf{b}$  are not related to  $\boldsymbol{\omega}$ , in this subsection, we investigate (20) instead of (19). The distributionally robust Bayesian estimation problem (20) can be rewritten as

$$\max_{\boldsymbol{\omega}} -\boldsymbol{\mu}^\top \mathbf{A}\boldsymbol{\mu} + \mathbf{b}^\top \boldsymbol{\mu} \quad (24)$$

$$\text{s.t.} \quad \begin{cases} \sum_{j=1}^N \omega_j = 1, \\ \omega_j \geq 0, \\ \Delta_0(\boldsymbol{\omega}, \bar{\boldsymbol{\omega}}) \leq \theta_0, \end{cases} \quad \forall j \in [N],$$

where  $\boldsymbol{\Theta}$  reduces to  $\boldsymbol{\omega}$ ; the relation between the prior model probability  $\omega_j$  and the posterior model probability  $\mu_j$ , for every model  $j$ , is established in (3);  $\mathbf{A}$  and  $\mathbf{b}$  are defined in (21) and (22), respectively, with  $\hat{\mathbf{x}}_j$  being replaced with  $\bar{\hat{\mathbf{x}}}_j$  and  $\mathbf{P}_j$  being replaced with  $\bar{\mathbf{P}}_j$  (cf. Appendix A). Note that in this case,  $\mathbb{P}_{j,\mathbf{x},\mathbf{y}} = \bar{\mathbb{P}}_{j,\mathbf{x},\mathbf{y}}, \forall j \in [N]$ .

Although the objective of (24) is quadratic in  $\boldsymbol{\mu}$  and both  $\mathbf{A}$  and  $\mathbf{b}$  are fixed, (24) is not a quadratic program because the decision vector is  $\boldsymbol{\omega}$  instead of  $\boldsymbol{\mu}$ , and the relation between  $\boldsymbol{\omega}$  and  $\boldsymbol{\mu}$  is highly nonlinear; see (3). Since uncertain prior model probabilities lead to uncertain posterior model probabilities, to simplify the problem, we propose to solve (24) over the posterior model probability  $\boldsymbol{\mu}$  rather than explicitly over the prior model probability  $\boldsymbol{\omega}$ .<sup>3</sup> We have

$$\max_{\boldsymbol{\mu}} -\boldsymbol{\mu}^\top \mathbf{A}\boldsymbol{\mu} + \mathbf{b}^\top \boldsymbol{\mu} \quad (25)$$

$$\text{s.t.} \quad \begin{cases} \sum_{j=1}^N \mu_j = 1 \\ \mu_j \geq 0, \\ \Delta_0(\boldsymbol{\mu}, \bar{\boldsymbol{\mu}}) \leq \theta_0. \end{cases} \quad \forall j \in [N],$$

Compared to the model (24), the model (25) is not only easier to be solved but also suitable to capture model uncertainties of the nominal model set because uncertain models give inexact model likelihood evaluations; cf. Step 1.6 in Algorithm 3 in the online supplementary materials. Since uncertainties in the posterior model probabilities may be caused by many factors (e.g., uncertain model transition probability matrix, uncertainties in candidate models), even though the state estimator for jump linear system (7) is robustified by taking into consideration only  $\boldsymbol{\mu}$ , the robust estimator is still likely to be insensitive to all the causing factors.

To measure the closeness between two discrete distributions  $\boldsymbol{\mu}$  and  $\bar{\boldsymbol{\mu}}$  (i.e., to quantify uncertainties in  $\boldsymbol{\mu}$ ), many statistical similarity measures can be used, for example, the popular Kullback–Leibler divergence and Wasserstein distance [30]; justifications and comparisons of using the two exemplified measures can be seen in, e.g., [31]–[36].

We use the Kullback–Leibler (KL) divergence to define  $\Delta_0(\cdot, \cdot)$ ; the case where the Wasserstein distance is employed to define  $\Delta_0(\cdot, \cdot)$  is discussed in Appendix C. The problem (25) is particularized to

$$\max_{\boldsymbol{\mu}} -\boldsymbol{\mu}^\top \mathbf{A}\boldsymbol{\mu} + \mathbf{b}^\top \boldsymbol{\mu} \quad (26)$$

$$\text{s.t.} \quad \begin{cases} \mathbf{1}^\top \boldsymbol{\mu} = 1 \\ \boldsymbol{\mu}^\top \ln \boldsymbol{\mu} - \boldsymbol{\mu}^\top \ln \bar{\boldsymbol{\mu}} \leq \theta_0, \end{cases}$$

where  $\ln \boldsymbol{\mu} := (\ln \mu_1, \ln \mu_2, \dots, \ln \mu_N)^\top$ . Since the objective is concave and the constraints are convex, the solution of (26) can be readily obtained using the Lagrangian duality theory (particularly, the Karush–Kuhn–Tucker conditions) [37]. For details, see Appendix I in the online supplementary materials. Note that (26) may have multiple optimal solutions. Note also that  $\boldsymbol{\mu} \geq \mathbf{0}$  is a redundant constraint because the function  $\ln \boldsymbol{\mu}$  implicitly requires it. We computationally prefer the KL divergence for  $\Delta_0(\cdot, \cdot)$  because the resulting solution method is computationally more efficient than that under the Wasserstein distance; for details, see Appendix C.

<sup>3</sup>For an extensive discussion, see Appendix H in the online supplementary materials.

B. When  $\theta_0 \neq 0$  And  $\theta_j \neq 0, \forall j \in [N]$

In this subsection, we discuss the case that both the nominal model set and the prior model probabilities are uncertain; i.e., the generic ambiguity set  $\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$  in (12) is considered.

*Proposition 3:* Suppose both the model set and the prior model probability  $\boldsymbol{\omega}$  are uncertain. If we conduct robustification over the posterior model probability  $\boldsymbol{\mu}$  instead of the prior model probability  $\boldsymbol{\omega}$ , the reformulated distributionally robust Bayesian estimation problem (19) can be written as

$$\begin{aligned} \max_{\boldsymbol{\Theta}'} \quad & \text{Tr} \sum_{j=1}^N \mu_j \cdot \{ \mathbf{P}_j + (\hat{\mathbf{x}}_j - \hat{\mathbf{x}})(\hat{\mathbf{x}}_j - \hat{\mathbf{x}})^\top \} \\ \text{s.t.} \quad & \begin{cases} \Delta_j(\mathbb{P}_{j,\mathbf{x},\mathbf{y}}, \bar{\mathbb{P}}_{j,\mathbf{x},\mathbf{y}}) \leq \theta_j, & \forall j \in [N] \\ \Delta_0(\boldsymbol{\mu}, \bar{\boldsymbol{\mu}}) \leq \theta_0, \quad \mathbf{1}^\top \boldsymbol{\mu} = 1, \quad \boldsymbol{\mu} \geq \mathbf{0}, \end{cases} \end{aligned} \quad (27)$$

where  $\hat{\mathbf{x}}$ ,  $\hat{\mathbf{x}}_j$ , and  $\mathbf{P}_j$  are defined in (17) and Appendix A, respectively;  $\boldsymbol{\Theta}' := \{\mu_j, \mathbf{c}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{x}}, \mathbf{c}_{j,\mathbf{v}}, \boldsymbol{\Sigma}_{j,\mathbf{v}}\}_{j \in [N]}$ . Recall that  $\mathbb{P}_{j,\mathbf{x},\mathbf{y}}$  is parameterized by  $\mathbf{c}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{x}}, \mathbf{c}_{j,\mathbf{v}}$ , and  $\boldsymbol{\Sigma}_{j,\mathbf{v}}$ , while  $\bar{\mathbb{P}}_{j,\mathbf{x},\mathbf{y}}$  is parameterized by  $\bar{\mathbf{x}}_j, \mathbf{M}_j$ , and  $\mathbf{R}_j$ ; see (16) and (15), respectively.

*Proof:* This is resulted from the definition of the ambiguity set  $\mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$  in (12).  $\square$

Note that  $\boldsymbol{\Theta}'$  in (27) is different from  $\boldsymbol{\Theta} := \{\omega_j, \mathbf{c}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{x}}, \mathbf{c}_{j,\mathbf{v}}, \boldsymbol{\Sigma}_{j,\mathbf{v}}\}_{j \in [N]}$  in (19). We optimize (27) over  $\boldsymbol{\Theta}'$  instead of  $\boldsymbol{\Theta}$  to reduce the complexity of the problem; however, one can also optimize (27) over  $\boldsymbol{\Theta}$  for the fidelity of the problem setting; see Appendix H in the online supplementary materials.

Due to the complexity of the definitions of  $\hat{\mathbf{x}}$  in (17) and  $\hat{\mathbf{x}}_j$  and  $\mathbf{P}_j$  in Appendix A, Problem (27) is difficult to solve; there exist strong nonlinearities and involve many decision variables. Hence, we simplify (27) by making the following assumption.

*Assumption 1:* Suppose  $\mathbf{c}_{j,\mathbf{x}} = \bar{\mathbf{x}}_j$ ,  $\mathbf{c}_{j,\mathbf{v}} = \mathbf{0}$ , and

$$\begin{aligned} \boldsymbol{\Sigma}_{j,\mathbf{x}} \mathbf{H}_j^\top (\mathbf{H}_j \boldsymbol{\Sigma}_{j,\mathbf{x}} \mathbf{H}_j^\top + \boldsymbol{\Sigma}_{j,\mathbf{v}})^{-1} \\ = \mathbf{M}_j \mathbf{H}_j^\top (\mathbf{H}_j \mathbf{M}_j \mathbf{H}_j^\top + \mathbf{R}_j)^{-1}, \end{aligned}$$

for every  $j \in [N]$ .  $\square$

Intuitively, when we conduct optimization in (27) by changing  $\{\mu_j, \mathbf{c}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{x}}, \mathbf{c}_{j,\mathbf{v}}, \boldsymbol{\Sigma}_{j,\mathbf{v}}\}_{j \in [N]}$ , Assumption 1 requires that the value of  $\hat{\mathbf{x}}_j$  remains unchanged and equal to its nominal value defined by  $\bar{\mathbf{x}}_j$  (cf. Appendix A) all the time. However, the value of  $\mathbf{P}_j$ , which is a function of variables  $\boldsymbol{\Sigma}_{j,\mathbf{x}}$  and  $\boldsymbol{\Sigma}_{j,\mathbf{v}}$ , would change. This is one of the practical yet promising tricks to shrink (i.e., limit the size of) the feasible region of (27).

By Assumption 1, the problem (27) can be written as

$$\begin{aligned} \max_{\boldsymbol{\Theta}'} \quad & \text{Tr} \sum_{j=1}^N \mu_j \cdot \{ \mathbf{P}_j + (\hat{\mathbf{x}}_j - \hat{\mathbf{x}})(\hat{\mathbf{x}}_j - \hat{\mathbf{x}})^\top \} \\ \text{s.t.} \quad & \begin{cases} \mathbf{c}_{j,\mathbf{x}} = \bar{\mathbf{x}}_j, \quad \mathbf{c}_{j,\mathbf{v}} = \mathbf{0}, & \forall j \in [N] \\ \boldsymbol{\Sigma}_{j,\mathbf{x}} \mathbf{H}_j^\top = \mathbf{K}_j \cdot (\mathbf{H}_j \boldsymbol{\Sigma}_{j,\mathbf{x}} \mathbf{H}_j^\top + \boldsymbol{\Sigma}_{j,\mathbf{v}}), & \forall j \in [N] \\ \Delta_j(\mathbb{P}_{j,\mathbf{x},\mathbf{y}}, \bar{\mathbb{P}}_{j,\mathbf{x},\mathbf{y}}) \leq \theta_j, & \forall j \in [N] \\ \Delta_0(\boldsymbol{\mu}, \bar{\boldsymbol{\mu}}) \leq \theta_0, \quad \mathbf{1}^\top \boldsymbol{\mu} = 1, \quad \boldsymbol{\mu} \geq \mathbf{0}, \end{cases} \end{aligned} \quad (28)$$

where

$$\mathbf{K}_j := \mathbf{M}_j \mathbf{H}_j^\top (\mathbf{H}_j \mathbf{M}_j \mathbf{H}_j^\top + \mathbf{R}_j)^{-1} \quad (29)$$

is a known matrix. As assumed, for every  $j \in [N]$ ,  $\mathbb{P}_{j,\mathbf{x},\mathbf{y}}$  and  $\bar{\mathbb{P}}_{j,\mathbf{x},\mathbf{y}}$  are Gaussian. Hence, the distance between  $\mathbb{P}_{j,\mathbf{x},\mathbf{y}}$

and  $\bar{\mathbb{P}}_{j,\mathbf{x},\mathbf{y}}$  is sufficient to be defined by the first two moments, i.e., their means and covariances. Recall that  $\mathbb{P}_{j,\mathbf{x},\mathbf{y}}$  is parameterized by  $\mathbf{c}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{x}}, \mathbf{c}_{j,\mathbf{v}}$ , and  $\boldsymbol{\Sigma}_{j,\mathbf{v}}$ , while  $\bar{\mathbb{P}}_{j,\mathbf{x},\mathbf{y}}$  is parameterized by  $\bar{\mathbf{x}}_j, \mathbf{M}_j$ , and  $\mathbf{R}_j$ . Therefore, it is sufficient to only define uncertainty sets for matrices  $\boldsymbol{\Sigma}_{j,\mathbf{x}}$  and  $\boldsymbol{\Sigma}_{j,\mathbf{v}}$  because  $\mathbf{c}_{j,\mathbf{x}}$  and  $\mathbf{c}_{j,\mathbf{v}}$  have been set to their nominal values. To be specific, the problem (28) can be rewritten as

$$\begin{aligned} \max_{\boldsymbol{\Theta}''} \quad & \text{Tr} \sum_{j=1}^N \mu_j \cdot \{ \mathbf{P}_j + (\hat{\mathbf{x}}_j - \hat{\mathbf{x}})(\hat{\mathbf{x}}_j - \hat{\mathbf{x}})^\top \} \\ \text{s.t.} \quad & \begin{cases} \boldsymbol{\Sigma}_{j,\mathbf{x}} \mathbf{H}_j^\top = \mathbf{K}_j \cdot (\mathbf{H}_j \boldsymbol{\Sigma}_{j,\mathbf{x}} \mathbf{H}_j^\top + \boldsymbol{\Sigma}_{j,\mathbf{v}}), & \forall j \in [N], \\ \Delta_{j,\mathbf{x}}(\boldsymbol{\Sigma}_{j,\mathbf{x}}, \mathbf{M}_j) \leq \theta_{j,\mathbf{x}}, & \forall j \in [N], \\ \Delta_{j,\mathbf{v}}(\boldsymbol{\Sigma}_{j,\mathbf{v}}, \mathbf{R}_j) \leq \theta_{j,\mathbf{v}}, & \forall j \in [N], \\ \Delta_0(\boldsymbol{\mu}, \bar{\boldsymbol{\mu}}) \leq \theta_0, \quad \mathbf{1}^\top \boldsymbol{\mu} = 1, \quad \boldsymbol{\mu} \geq \mathbf{0}, \\ \boldsymbol{\Sigma}_{j,\mathbf{x}} \succeq \mathbf{0}, \boldsymbol{\Sigma}_{j,\mathbf{v}} \succeq \mathbf{0}, & \forall j \in [N], \end{cases} \end{aligned} \quad (30)$$

where  $\boldsymbol{\Theta}'' := \{\mu_j, \boldsymbol{\Sigma}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{v}}\}_{j \in [N]}$ ,  $\Delta_{j,\mathbf{x}}(\boldsymbol{\Sigma}_{j,\mathbf{x}}, \mathbf{M}_j)$  means a similarity measure between the two matrices  $\boldsymbol{\Sigma}_{j,\mathbf{x}}$  and  $\mathbf{M}_j$ , and  $\Delta_{j,\mathbf{v}}(\boldsymbol{\Sigma}_{j,\mathbf{v}}, \mathbf{R}_j)$  denotes a similarity measure between the two matrices  $\boldsymbol{\Sigma}_{j,\mathbf{v}}$  and  $\mathbf{R}_j$ .

The problem (30) can be solved using the coordinate descent method: we can first fix  $\boldsymbol{\mu}$  and optimize over  $\{\boldsymbol{\Sigma}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{v}}\}_{j \in [N]}$ , and then fix  $\{\boldsymbol{\Sigma}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{v}}\}_{j \in [N]}$  and optimize over  $\boldsymbol{\mu}$ . The iterative process proceeds until it converges. Since  $\mu_j > 0, \forall j \in [N]$  and the constraints of  $\boldsymbol{\mu}$  are not coupled with the constraints of  $\{\boldsymbol{\Sigma}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{v}}\}_{j \in [N]}$ , this iterative process can converge only in one round. To be specific, it is sufficient to first fix  $\boldsymbol{\mu}$  and optimize over  $\{\boldsymbol{\Sigma}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{v}}\}_{j \in [N]}$ , and then fix  $\{\boldsymbol{\Sigma}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{v}}\}_{j \in [N]}$  and optimize over  $\boldsymbol{\mu}$ : no more iterations are required.<sup>4</sup>

The sub-problem where  $\{\boldsymbol{\Sigma}_{j,\mathbf{x}}\}_{j \in [N]}$  and  $\{\boldsymbol{\Sigma}_{j,\mathbf{v}}\}_{j \in [N]}$  are fixed and the optimization is conducted over  $\boldsymbol{\mu}$  has been solved in Subsection IV-A. In the following, we focus on solving the sub-problem when  $\boldsymbol{\mu}$  is fixed, which is

$$\begin{aligned} \max_{\boldsymbol{\Upsilon}} \quad & \sum_{j=1}^N \mu_j \cdot \text{Tr} \mathbf{P}_j \\ \text{s.t.} \quad & \begin{cases} \boldsymbol{\Sigma}_{j,\mathbf{x}} \mathbf{H}_j^\top = \mathbf{K}_j \cdot (\mathbf{H}_j \boldsymbol{\Sigma}_{j,\mathbf{x}} \mathbf{H}_j^\top + \boldsymbol{\Sigma}_{j,\mathbf{v}}), & \forall j \in [N], \\ \Delta_{j,\mathbf{x}}(\boldsymbol{\Sigma}_{j,\mathbf{x}}, \mathbf{M}_j) \leq \theta_{j,\mathbf{x}}, \quad \boldsymbol{\Sigma}_{j,\mathbf{x}} \succeq \mathbf{0}, & \forall j \in [N], \\ \Delta_{j,\mathbf{v}}(\boldsymbol{\Sigma}_{j,\mathbf{v}}, \mathbf{R}_j) \leq \theta_{j,\mathbf{v}}, \quad \boldsymbol{\Sigma}_{j,\mathbf{v}} \succeq \mathbf{0}, & \forall j \in [N], \end{cases} \end{aligned} \quad (31)$$

where  $\boldsymbol{\Upsilon} := \{\boldsymbol{\Sigma}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{v}}\}_{j \in [N]}$ . Since  $\boldsymbol{\mu} > \mathbf{0}$ , for each  $j \in [N]$ , the problem (31) can be separately solved because both the objective function and the feasible region are separable. Specifically, for all  $i, j \in [N]$  and  $i \neq j$ ,  $\mathbf{P}_i$  and  $\mathbf{P}_j$  are not coupled, so are their respective constraints. As a result, the solution  $\boldsymbol{\Upsilon}^* = \{\boldsymbol{\Sigma}_{j,\mathbf{x}}^*, \boldsymbol{\Sigma}_{j,\mathbf{v}}^*\}_{j \in [N]}$  of (31) also solves

$$\begin{aligned} \max_{\boldsymbol{\Sigma}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{v}}} \quad & \text{Tr} \mathbf{P}_j \\ \text{s.t.} \quad & \begin{cases} \boldsymbol{\Sigma}_{j,\mathbf{x}} \mathbf{H}_j^\top = \mathbf{K}_j \cdot (\mathbf{H}_j \boldsymbol{\Sigma}_{j,\mathbf{x}} \mathbf{H}_j^\top + \boldsymbol{\Sigma}_{j,\mathbf{v}}), \\ \Delta_{j,\mathbf{x}}(\boldsymbol{\Sigma}_{j,\mathbf{x}}, \mathbf{M}_j) \leq \theta_{j,\mathbf{x}}, \quad \boldsymbol{\Sigma}_{j,\mathbf{x}} \succeq \mathbf{0} \\ \Delta_{j,\mathbf{v}}(\boldsymbol{\Sigma}_{j,\mathbf{v}}, \mathbf{R}_j) \leq \theta_{j,\mathbf{v}}, \quad \boldsymbol{\Sigma}_{j,\mathbf{v}} \succeq \mathbf{0} \end{cases} \end{aligned} \quad (32)$$

<sup>4</sup>This can be intuitively seen from the equivalence between (31) and (32): the specific value of  $\boldsymbol{\mu}$  does not influence the optimal values of  $\{\boldsymbol{\Sigma}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{v}}\}_{j \in [N]}$ . Hence, the optimization (30) over  $\{\boldsymbol{\Sigma}_{j,\mathbf{x}}, \boldsymbol{\Sigma}_{j,\mathbf{v}}\}_{j \in [N]}$  can be completely solved in the first-round iteration. Afterward, the optimization over  $\boldsymbol{\mu}$  can be readily and completely solved as well.

for every  $j \in [N]$ , and vice versa. Therefore, it necessitates and suffices to separately solve (32) for every  $j \in [N]$ .

Recall from Appendix A for the definition of the matrix function:

$$\begin{aligned} P_j(\Sigma_{j,x}, \Sigma_{j,v}) \\ := \Sigma_{j,x} - \Sigma_{j,x} H_j^\top (H_j \Sigma_{j,x} H_j^\top + \Sigma_{j,v})^{-1} H_j \Sigma_{j,x}. \end{aligned}$$

The theorem below plays an important role to further simplify the problem (32).

*Theorem 1:* The matrix function  $\text{Tr } P_j(\Sigma_{j,x}, \Sigma_{j,v})$  is monotonically increasing in both  $\Sigma_{j,x}$  and  $\Sigma_{j,v}$ . Namely, for every given  $\Sigma_{j,v}$ , if  $\Sigma_1 \succeq \Sigma_2$ , we have  $\text{Tr } P_j(\Sigma_1, \Sigma_{j,v}) \geq \text{Tr } P_j(\Sigma_2, \Sigma_{j,v})$ . Likewise, for every given  $\Sigma_{j,x}$ , if  $\Sigma_1 \succeq \Sigma_2$ , we have  $\text{Tr } P_j(\Sigma_{j,x}, \Sigma_1) \geq \text{Tr } P_j(\Sigma_{j,x}, \Sigma_2)$ .

*Proof:* See Appendix E.  $\square$

Theorem 1 equivalently transforms the problem (32) to

$$\begin{aligned} \max_{\Sigma_{j,x}, \Sigma_{j,v}} \quad & \text{Tr}[\Sigma_{j,x} + \Sigma_{j,v}] \\ \text{s.t.} \quad & \begin{cases} \Sigma_{j,x} H_j^\top = K_j \cdot (H_j \Sigma_{j,x} H_j^\top + \Sigma_{j,v}), \\ \Delta_{j,x}(\Sigma_{j,x}, M_j) \leq \theta_{j,x}, \quad \Sigma_{j,x} \succeq 0, \\ \Delta_{j,v}(\Sigma_{j,v}, R_j) \leq \theta_{j,v}, \quad \Sigma_{j,v} \succeq 0. \end{cases} \end{aligned} \quad (33)$$

The proposition below summarizes the solution of the reformulated distributionally robust Bayesian estimation problem (27).

*Proposition 4:* Suppose both the model set and the prior model probability  $\omega$  are uncertain. If we conduct the robustification over the posterior model probability  $\mu$  instead of the prior model probability  $\omega$  and Assumption 1 is adopted, then the reformulated distributionally robust Bayesian estimation problem (27) can be solved by

- 1) **Step 1.** Solving (33) for every  $j \in [N]$ ;
- 2) **Step 2.** Solving (25) with updated  $\{\Sigma_{j,x}, \Sigma_{j,v}\}_{\forall j \in [N]}$  from Step 1.<sup>5</sup>  $\square$

In the following, as a demonstration, we discuss the particular cases where  $\Delta_{j,x}(\cdot, \cdot)$  and  $\Delta_{j,v}(\cdot, \cdot)$  are defined by the moments-based similarity measure. The case where  $\Delta_{j,x}(\cdot, \cdot)$  and  $\Delta_{j,v}(\cdot, \cdot)$  are defined by the Wasserstein distance is discussed in Appendix F.

If we use the moment-based ambiguity sets for  $\mathbb{P}_{j,x,y}$  as in [13], the problem (33) is particularized into

$$\begin{aligned} \max_{\Sigma_{j,x}, \Sigma_{j,v}} \quad & \text{Tr}[\Sigma_{j,x} + \Sigma_{j,v}] \\ \text{s.t.} \quad & \begin{cases} \Sigma_{j,x} H_j^\top = K_j \cdot (H_j \Sigma_{j,x} H_j^\top + \Sigma_{j,v}), \\ (1 - \theta_{j,x}) M_j \preceq \Sigma_{j,x} \preceq (1 + \theta_{j,x}) M_j, \\ (1 - \theta_{j,v}) R_j \preceq \Sigma_{j,v} \preceq (1 + \theta_{j,v}) R_j, \\ \Sigma_{j,x} \succeq 0, \quad \Sigma_{j,v} \succeq 0, \end{cases} \end{aligned} \quad (34)$$

where  $0 \leq \theta_{j,x} \leq 1$ ,  $0 \leq \theta_{j,v} \leq 1$ . The problem (34) is a standard linear positive semi-definite program (SDP), which can be efficiently solved using mature SDP solvers such as MOSEK, GUROBI, YALMIP. However, we show that in a special case, (34) can be analytically solved.

*Theorem 2:* If  $\theta_{j,x} = \theta_{j,v} = \theta_j$ , the optimal solution to (34) is analytically given by

$$\begin{cases} \Sigma_{j,x}^* = (1 + \theta_j) M_j, \\ \Sigma_{j,v}^* = (1 + \theta_j) R_j. \end{cases} \quad (35)$$

<sup>5</sup>When  $\{\Sigma_{j,x}, \Sigma_{j,v}\}_{\forall j \in [N]}$  are updated,  $A$  and  $b$  in (25) would be updated as well because  $\hat{x}_j, P_j, \forall j \in [N]$  therein would be updated.

*Proof:* See Appendix G.  $\square$

## V. DISTRIBUTIONALLY ROBUST INTERACTIVE MULTIPLE MODEL STATE ESTIMATOR

The overall distributionally robust interactive multiple model (DRIMM) state estimator is summarized in Algorithm 1, which is a robustified, modified version of the standard IMM filter in Algorithm 3 in the online supplementary materials. In Algorithm 1, for clarity, we only highlight the modifications and differences between the proposed distributionally robust IMM filter and the standard IMM filter in Algorithm 3. For other algorithmic statements, see Algorithm 3. Recall from Appendix A that  $P_j$  is defined by  $\Sigma_{j,x}$  and  $\Sigma_{j,v}$ . In Step 2.1, another modification method for  $P_{j,k|k}$  based on Wasserstein distance is available in Appendix F; in Step 2.2, another modification method for  $\mu_{k|k}$  based on Wasserstein distance is available in Appendix C. Note that Step 2.1 in Algorithm 1 solves (34) using (35) to modify  $P_{j,k|k}, \forall j \in [N]$ ; Line 6 is due to Theorem 2 and the definition of  $P_{j,k|k}$  is in Appendix A. Note also that the solution of (26) is given in Algorithm 2 in the online supplementary materials.

### Algorithm 1 Distributionally Robust IMM State Estimator

**Definition:**  $\theta_0, \theta_j$ , for every  $j \in [N]$  are the size parameters to define the scale of the ambiguity set (12).  $N$  is the number of the nominal models.  $k$  denotes the discrete time index.  $P_{j,k|k}$  is the posterior state estimation error covariance of the  $j^{\text{th}}$  nominal model at the time  $k$ .  $\mu_{k|k}$  is the posterior model probability at the time  $k$ .

**External:** Algorithm 3 (i.e., the standard IMM filter [3], [5]) in the online supplementary materials.

**Input:**  $\theta_0, \theta_j$ , for every  $j \in [N]$ .

```

1: while true do
2:   // (Step 1) At Time  $k$ 
3:   Execute Step 1 (i.e., Lines 2 ~ 24) of Algorithm 3
4:   // (Step 2) Robustification; See Proposition 4
5:   // (Step 2.1) Modify  $P_{j,k|k}, \forall j \in [N]$  via (34)
6:    $P_{j,k|k} \leftarrow (1 + \theta_j) \cdot P_{j,k|k}, \forall j \in [N]$ 
7:   // (Step 2.2) Modify  $\mu_{k|k}$  via (25)
8:   Solve (26) to modify  $\mu_{k|k}$ 
9:   // (Step 3) Combined Posterior State Estimate
10:  Execute Step 2 (i.e., Lines 25 ~ 27) of Algorithm
    3 to calculate robustified  $\hat{x}_{k|k}$  and  $P_{k|k}$ , using modified
     $P_{j,k|k}, \forall j \in [N]$  and  $\mu_{k|k}$ 
11:  // Next Time Step
12:   $k \leftarrow k + 1$ 
13: end while
```

**Output:**  $\hat{x}_{k|k}, P_{k|k}, \mu_{k|k}, \forall k$

*Remark 2 (Only TPMs Are Uncertain):* Uncertain TPMs lead to uncertain model weights. If there do not exist uncertainties in candidate models and only model weights are uncertain, we can robustify the IMM filter through only modifying posterior model weights  $\mu_{k|k}$ . In this case, Step 2.1 in Algorithm 1 should be ignored and only Step 2.2 should be conducted.  $\square$



*Remark 3 (Candidate Models Are Uncertain):* As long as there exist uncertainties in the designed nominal model set, the DRIMM filter in Algorithm 1 cannot be simplified: Neither Step 2.1 nor Step 2.2 can be ignored. This is because uncertain candidate models introduce uncertainties in evaluating posterior model weights. Specifically, from (3), we can see that when there exist uncertainties in candidate models, the evaluation of model likelihoods cannot be exact, and therefore, the posterior model weights cannot be exact either. However, in practice, when candidate models are uncertain, one may ignore Step 2.2 and just conduct Step 2.1. This leads to a simplification (or approximation) of Algorithm 1, which is an IMM filter using the robustified Kalman filter (rather than the conventional non-robust one) for each candidate model. Recall that the distributionally robust state estimation method for each candidate model has been discussed in [13], [14].  $\square$

*Remark 4 (Computational Burdens):* Compared with the standard IMM filter in Algorithm 3 in the online supplementary materials, the DRIMM filter in Algorithm 1 has an additional robustification procedure in Step 2. In Step 2.1, there is only a value assignment procedure, and hence, no extra computational burdens are introduced. However, in Step 2.2, an iterative process is involved to numerically solve (26); see Algorithm 2 in the online supplementary materials. Therefore, extra computational resources are required to conduct robustification over  $\mu_{k|k}$ .  $\square$

## VI. EXPERIMENTS

We conduct multi-model target tracking experiments to show the power of the proposed method. All the source data and codes are available online at GitHub: <https://github.com/Sprtm-Asleaf/DRSE-Jump>. The results are obtained by a Lenovo laptop with 16G RAM and 11th Gen Intel(R) Core(TM) i5-11300H CPU @ 3.10GHz. The programming environment is MATLAB 2019B. We implement the following methods for performance comparison.

- 1) **KF**: the standard Kalman filter with the exactly known model transition trajectory. In practice, the standard Kalman filter is *not applicable* for jump linear systems because we do not know the true model transition trajectory. However, in a simulation experiment, we know the underlying true model transition trajectory so that the standard Kalman filter can be used to provide the theoretically optimal filtering performance.
- 2) **IMM-T**: the IMM filter using the exactly *true* model transition probability matrix. Note that in a simulation experiment, we know the underlying true model transition probability matrix. But in practice, this method is *not applicable* as the standard Kalman filter is.
- 3) **IMM-N**: the IMM filter using a user-specified *nominal* model transition probability matrix.
- 4) **IMM-R**: the IMM filter using the distributionally *robust* Kalman filter in [13] for every operating mode  $j \in [N]$ . To clarify further, in Algorithm 3 in supplementary materials, Step 1.4 is modified according to the robustified Kalman filter in [13] and all other steps in Algorithm 3 remain unchanged. See also Remark 3.

- 5) **IMM-B**: the IMM filter using the *Bayesian* method to estimate the true TPM [7]. This method is applicable only when TPM is uncertain.
- 6) **IMM-M**: the IMM filter using the *maximum* likelihood method to estimate the true TPM [12]. This method is applicable only when TPM is uncertain.
- 7) **DRIMM**: the distributionally robust IMM filter in Algorithm 1, which shares the same nominal model transition probability matrix with the IMM-N filter. This method is applicable for any types of uncertainty. Particularly, it *can* flexibly respond to specified types of uncertainty.
- 8) **IMM-RS**: the *risk-sensitive* IMM filter [21]. This method is applicable for any types of uncertainty. However, it *cannot* flexibly respond to specified types of uncertainty.
- 9) **IMM-C**: the *compensation*-based IMM filter [9]. This method is applicable for any types of uncertainty. However, it *cannot* flexibly respond to specified types of uncertainty.

Since the KF and the IMM-T filters are *not applicable* in practice, given a nominal TPM, a method is promising if it can outperform the IMM-N filter. Note that the IMM-T filter is *not optimal* for jump linear systems due to its approximation nature as elucidated in the introduction.

### A. Simulated Experiments

We continue studying the classic one-dimensional target tracking problem in [7], [9], [11], [17], [38] where the target maneuvers with Markov switching accelerations. According to [2], the motion of a target can be independently tracked in each axis; therefore, focusing only on one axis does not lose the generality. The jump Markov linear system is defined by

$$\begin{cases} \begin{bmatrix} p_{k+1} \\ s_{k+1} \end{bmatrix} = \begin{bmatrix} 1 & T \\ 0 & 1 \end{bmatrix} \begin{bmatrix} p_k \\ s_k \end{bmatrix} + \begin{bmatrix} T^2/2 \\ T \end{bmatrix} [a_{j,k} + w_k] \\ y_k = p_k + v_k \end{cases} \quad (36)$$

where  $p_k \in \mathbb{R}$  and  $s_k \in \mathbb{R}$  denote the position and velocity of a moving target at time  $k$ , respectively;  $T$  denotes the sampling time;  $a_{j,k} \in \mathbb{R}$  denotes the possible target maneuvering acceleration at time  $k$  which takes the value of

$$a_{j,k} = \begin{cases} 0, & j = 1, \\ 20, & j = 2, \\ -20, & j = 3; \end{cases} \quad (37)$$

$w_k \in \mathbb{R}$  and  $v_k \in \mathbb{R}$  denote the acceleration noise and the measurement noise at the time  $k$ , respectively. The following settings, as conventionally made in the mentioned state-of-the-art, are taken in this section:  $p_0 \sim \mathcal{N}(80000, 100^2)$ ,  $s_0 \sim \mathcal{N}(400, 100^2)$ ,  $w_k \sim \mathcal{N}(0, 2^2)$ ,  $v_k \sim \mathcal{N}(0, 100^2)$ , for every  $k$ , and the model's initial probability  $\mu_0 = [0.8, 0.1, 0.1]^T$ . In this paper, we set  $T = 1$ . Throughout the simulated experiments, we use  $\Pi_0$  to denote the true model transition probability matrix (TPM) and  $\Pi$  the nominal TPM. In addition, the nominal value of maneuvering acceleration is denoted as  $\bar{a}_{j,k}$ , which might be different from the true  $a_{j,k}$  in (37).

For each simulation scenario, we conduct 100 independent Monte-Carlo episodes and each episode runs 1000 time steps.

The overall performance evaluation method for each filter is the averaged root-mean-square error (RMSE):

$$\frac{1}{100} \sum_{l=1}^{100} \left[ \sqrt{\frac{1}{1000} \sum_{k=1}^{1000} (p_k^{(l)} - \hat{p}_k^{(l)})^2 + (s_k^{(l)} - \hat{s}_k^{(l)})^2} \right],$$

where  $p_k^{(l)}$  (resp.  $s_k^{(l)}$ ) denotes the true value of position (resp. velocity) at the time  $k$  in the  $l^{\text{th}}$  Monte-Carlo simulation, and  $\hat{p}_k^{(l)}$  (resp.  $\hat{s}_k^{(l)}$ ) its estimate. Every method is well-tuned and performs at its best for the given simulation scenario; details can be accessed in the shared source codes and therefore omitted here. For example, in the proposed DRIMM filter (Algorithm 1), we have set  $\theta_0 = 0.1$  and  $\theta_j = 0.25$  if the  $j^{\text{th}}$  candidate model is uncertain. The influence of the values of  $\theta$ 's to the filtering performances will be discussed later in Subsection VI-A5.

1) *When Only TPM Is Uncertain:* First, we suppose that the nominal  $\Pi$  is different from the true  $\Pi_0$ . However, the candidate models are exact (i.e., the nominal  $\bar{a}_{j,k}$  is the same as the true  $a_{j,k}$  in (37)). In this case, the IMM-R method is not applicable since it is not designed to handle the uncertainty in  $\Pi$ . Note that Step 2.1 in Algorithm 1 should be ignored; see Remark 2.

We first set the true  $\Pi_0$  and the nominal  $\Pi$  as follows:

$$\Pi_0 := \begin{bmatrix} 0.8 & 0.1 & 0.1 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}, \Pi := \begin{bmatrix} 0.9 & 0.05 & 0.05 \\ 0.05 & 0.9 & 0.05 \\ 0.05 & 0.05 & 0.9 \end{bmatrix}.$$

The results are shown in Table I.

TABLE I  
WHEN ONLY TPM IS UNCERTAIN

Filter	RMSE	Avg Time	Filter	RMSE	Avg Time
KF*	43.62	1.77e-06	IMM-T*	81.60	1.99e-05
IMM-N	82.44	1.94e-05	IMM-B	82.36	2.54e-05
IMM-M	82.32	5.32e-05	IMM-C	166.01	2.06e-05
IMM-RS	82.37	8.05e-05	DRIMM	<b>82.19</b>	2.14e-04

**Avg Time:** Average Execution Time at each time step (unit: seconds);  
**1e-5:**  $1 \times 10^{-5}$ ;

\*: Not applicable methods in practice.

Next, we set  $\Pi_0$  and  $\Pi$  as follows:

$$\Pi_0 := \begin{bmatrix} 0.1 & 0.1 & 0.8 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}, \Pi := \begin{bmatrix} 0.8 & 0.1 & 0.1 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}.$$

That is, the first row and the third row of  $\Pi_0$  are the same. However, the nominal  $\Pi$  mistakenly assumes that they are different. The results are given in Table II.

TABLE II  
WHEN ONLY TPM IS UNCERTAIN (ANOTHER EXAMPLE)

Filter	RMSE	Avg Time	Filter	RMSE	Avg Time
KF*	43.62	1.91e-06	IMM-T*	82.48	2.00e-05
IMM-N	84.56	1.98e-05	IMM-B	84.03	2.54e-05
IMM-M	<b>83.95</b>	5.41e-05	IMM-C	111.60	2.07e-05
IMM-RS	84.63	8.09e-05	DRIMM	84.47	2.91e-04

See Table I for table notes.

From Tables I and II (and also many other similar experiments), the following main points have to be highlighted.

- 1) The compensation-based method (i.e., IMM-C) is not satisfactory for the two cases. This is because there does not exist a dominating operating mode for both two cases: i.e., the systems' true operating modes are frequently switching. (However, when the true switching frequency is low, the IMM-C method would be useful; see the real-world target tracking experiments in Appendix VI-B.)
- 2) In the first case (Table I), the IMM-B method, the IMM-M method, the IMM-RS method, and the proposed DRIMM method are all better than the nominal IMM-N method. This suggests that the four methods are all able to combat uncertainties in the nominal TPM  $\Pi$ . In the second case (Table II), the IMM-M and the IMM-B methods perform significantly better than the IMM-RS, the proposed DRIMM, and the nominal IMM-N methods. This is because the IMM-M and the IMM-B methods can adaptively estimate the unknown true TPM so that the uncertainties in the nominal TPM would be reduced. As a result, the filtering performances of the IMM-M and the IMM-B methods can be significantly improved. This result is consistent with the findings in [13], [30]: i.e., adaptive methods which aim to reduce the uncertainties would be potentially better than robust methods that aim to just hedge against uncertainties.
- 3) Among the robust methods, the proposed DRIMM method performs better than the generic IMM-RS method because the former is able to specifically respond to the uncertainties only in the nominal TPM. In contrast, the IMM-RS method only assumes that there exist model uncertainties; it does not take into account where and how the uncertainties exist. Therefore, the IMM-RS filter tends to be overly conservative.

In summary, when there exist uncertainties only in the nominal TPM, we would suggest that the practitioners should first consider adaptive methods such as the IMM-M and the IMM-B methods because the two methods tend to estimate the unknown true TPM and the uncertainties in the nominal TPM are to be reduced. In addition, the two adaptive methods are computationally more efficient compared to the proposed DRIMM method.

2) *When The Model Set Is Incomplete:* Second, we consider the case where the true system dynamics has five operating modes:  $a_{1,k} = 0$ ,  $a_{2,k} = 10$ ,  $a_{3,k} = -10$ ,  $a_{4,k} = 20$ , and  $a_{5,k} = -20$ . But the nominal model set still assumes that there exist three operating modes, as in (37). This is a common situation in target tracking. To be specific, when we do not know the true maneuvering acceleration  $a_k$  of the moving target, there exist infinitely many possible values for  $a_k$ . However, considering practical constraints such as computational complexity, we may only assume finitely many values for  $a_k$ . As a result, the designed nominal model set is potentially incomplete. The true and nominal TPMs are set, respectively, to

$$\Pi_0 := \begin{bmatrix} 0.6 & 0.1 & 0.1 & 0.1 & 0.1 \\ 0.1 & 0.6 & 0.1 & 0.1 & 0.1 \\ 0.1 & 0.1 & 0.6 & 0.1 & 0.1 \\ 0.1 & 0.1 & 0.1 & 0.6 & 0.1 \\ 0.1 & 0.1 & 0.1 & 0.1 & 0.6 \end{bmatrix}, \quad \Pi := \begin{bmatrix} 0.6 & 0.2 & 0.2 \\ 0.2 & 0.6 & 0.2 \\ 0.2 & 0.2 & 0.6 \end{bmatrix}.$$

The results are shown in Table III.

TABLE III  
THE NOMINAL MODEL SET IS NOT COMPLETE

Filter	RMSE	Avg Time	Filter	RMSE	Avg Time
KF*	43.62	1.77e-06	IMM-T*	92.24	3.79e-05
IMM-N	101.04	1.94e-05	IMM-B	100.73	2.54e-05
IMM-M	100.00	5.49e-05	IMM-C	100.42	2.09e-05
IMM-RS	98.53	8.02e-05	DRIMM	<b>95.18</b>	2.98e-04

See Table I for table notes.

As we can see from Table III, in this case, the adaptive methods (i.e., the IMM-B and IMM-M methods) cannot work significantly better than the nominal IMM-N method because if the nominal model set is incomplete, nominal candidate models are inexact models when the true operating modes are not included in the nominal model set. Hence, there exist uncertainties not only in the TPM but also in nominal candidate models. As a result, robust methods (i.e., the IMM-RS and DRIMM methods) that are able to combat all kinds of uncertainties are promising. Among robust methods, the proposed DRIMM method outperforms the IMM-RS method.

3) *When Candidate Models Are Uncertain:* Third, we suppose that the candidate models are uncertain (i.e., the nominal  $\bar{a}_{j,k} = [0, 10, -10]$  is different from the true  $a_{j,k} = [0, 20, -20]$  in (37)).

**When  $\Pi = \Pi_0$ .** Let the diagonal elements of the TPM be 0.8s, and all other non-diagonal elements be 0.1s. In this case, the IMM-B method and the IMM-M are not applicable since they are not designed to handle the uncertainties in candidate models. The results are shown in Table IV.

TABLE IV  
CANDIDATE MODELS ARE UNCERTAIN BUT TPM IS EXACT

Filter	RMSE	Avg Time	Filter	RMSE	Avg Time
KF*	43.62	1.80e-06	IMM-T*	81.60	2.02e-05
IMM-N	95.83	2.08e-05	IMM-C	161.48	2.15e-05
IMM-RS	90.43	8.44e-05	DRIMM	<b>87.67</b>	2.85e-04
IMM-R	88.14	2.54e-05			

See Table I for table notes.

As we can see from Table IV, when there exist uncertainties only in the candidate models, the IMM-R filter, the IMM-RS filter, and the DRIMM filter can outperform the nominal IMM-N filter. In addition, the proposed DRIMM filter works better than the IMM-R filter, which validates the claims in Remark 3: That is, uncertain candidate models lead to uncertain model weights, and therefore, robustification should also be conducted over model weights.

**When  $\Pi \neq \Pi_0$ .** Let the diagonal elements of the true TPM (i.e.,  $\Pi_0$ ) be 0.8s, and all other non-diagonal elements be 0.1s. Let the diagonal elements of the nominal TPM (i.e.,  $\Pi$ ) be 0.6s, and all other non-diagonal elements be 0.2s. The results are shown in Table V.

TABLE V  
BOTH CANDIDATE MODELS AND TPM ARE UNCERTAIN

Filter	RMSE	Avg Time	Filter	RMSE	Avg Time
KF*	43.62	1.84e-06	IMM-T*	81.60	1.98e-05
IMM-N	104.48	1.92e-05	IMM-B	103.99	2.61e-05
IMM-M	102.93	5.35e-05	IMM-C	103.15	2.07e-05
IMM-RS	98.44	7.98e-05	DRIMM	<b>90.23</b>	2.99e-04
IMM-R	93.25	2.50e-05			

See Table I for table notes.

As we can see from Table V, when there exist uncertainties in both the candidate models and the nominal TPM, all methods (except the KF and IMM-T methods) can outperform the nominal IMM-N filter because the uncertainties can be adaptively reduced or robustly withstand to some degree. Also, the proposed DRIMM method works better than the IMM-RS method as the latter tends to be overly conservative. In addition, the DRIMM method is better than the IMM-R method because the latter can only partially respond to the uncertainties in the candidate models. In contrast, the DRIMM method is able to work for uncertainties in both the candidate models and the nominal TPM.

4) *When There Are No Model Uncertainties:* Fourth, we suppose that there are no any model uncertainties (i.e., both the TPM and all candidate models are exact). We let the diagonal elements of the TPM (N.B.:  $\Pi_0 = \Pi$ ) be 0.8s, and all other non-diagonal elements be 0.1s. In addition, we let the nominal  $\bar{a}_k$  be equal to the true  $a_k$  in (37). The filtering results are shown in Table VI.

TABLE VI  
THERE EXIST NO ANY MODEL UNCERTAINTIES

Filter	RMSE	Avg Time	Filter	RMSE	Avg Time
KF*	43.62	1.78e-06	IMM-T*	81.60	2.01e-05
IMM-N	<b>81.60</b>	1.98e-05	IMM-B	81.61	2.59e-05
IMM-M	81.61	5.37e-05	IMM-C	110.69	2.18e-05
IMM-RS	82.39	8.05e-05	DRIMM	82.45	1.67e-04
IMM-R	82.44	2.51e-05			

See Table I for table notes.

As we can see from Table VI, when there exist no model uncertainties, all modified IMM methods would perform worse than the nominal IMM-N method. This result is consistent with the findings in [13], [30], [39]: There exists a trade-off between the robustness in uncertain conditions and the optimality in perfect conditions. To be specific, robust methods are not optimal in exact conditions but they can withstand uncertainties in inexact conditions; in contrast, the nominal IMM-N method is the best in exact conditions but it has no ability to combat uncertainties in inexact conditions.

5) *On The Size of Uncertainty Set (12):* In this subsection, we investigate the influence of the values of the size parameters  $\theta$ 's in the uncertainty set (12). As an example and without loss of generality, we suppose that the TPM is exact and only the candidate models are uncertain. To be specific, the nominal  $\bar{a}_{j,k} = [0, 10, -10]$  is different from the true  $a_{j,k} = [0, 20, -20]$ . The diagonal elements of the TPM (N.B.:  $\Pi_0 = \Pi$ ) are set to 0.8s and the non-diagonal ones

are set to 0.1s. The performances of the DRIMM filter are given in Table VII, against the value of  $\theta := \theta_j$  for every  $j$  such that the  $j$ th model is uncertain (i.e.,  $j = 2, 3$  in this case). As we can see from Table VII, the size of the ambiguity set (12) can be neither too large nor too small. If the size is too large, the DRIMM filter would be extremely conservative, while if the size is too small, the DRIMM filter cannot offer sufficient robustness. This result is consistent with those in [13], [30], [39]. However, the convincing tuning method for  $\theta$ 's is lacking because the true system state (i.e., training data set) is unknown for a real state estimation problem, and therefore,  $\theta$ 's cannot be rigorously tuned to be optimal. As such, we leave this as an open problem to be solved by the future research. At present, we can only suggest that practitioners can try appropriate values of  $\theta$ 's for their specific problems.

TABLE VII  
PERFORMANCES (I.E., RMSE) OF DRIMM AGAINST VALUES OF  $\theta$

$\theta$	0.0	0.1	0.2	0.3	0.4
RMSE	95.44	90.49	87.67	86.10	<b>85.30</b>
$\theta$	0.5	0.6	0.7	0.8	0.9
RMSE	85.41	85.36	85.51	85.82	86.24

### B. Real-World Target Tracking Examples

We track the real-time positions and velocities of a slowly-maneuvering car and a highly-maneuvering quadrotor drone, respectively. The raw position measurements are received by traditional GPS (global positioning system) solutions, while the high-accuracy real-time position and velocity measurements (treated as ground truths) are obtained by RTK (real-time kinematic) solutions [40]. We continue adopting the target tracking framework in (36) where the target's true acceleration at every time  $k$  is unknown, and hence, some nominal values are used. The nominal model set is therefore incomplete and uncertain.

1) *Track A Slowly-Maneuvering Car*: In this experiment, the data are collected by a slowly-maneuvering car that carries a GPS solution and an RTK solution; the car travels on a road in Beijing, China. The commercial model of the used RTK chip is P327 and of the antenna is UA35, both produced by UniStrong Co., Ltd., Beijing, China; see <http://en.unistrong.com/>. The car, its real-time velocities, and its trajectory are shown in Fig. 1.

Without loss of generality, we study the tracking problem in the east axis (in the east-north-up coordinate). We suppose that there are three nominal values for acceleration in the east axis:  $\bar{a}_{j,k} = [0, 2.5, -2.5]$ . The diagonal elements of the nominal TPM are set to 0.8s and the non-diagonal ones are set to 0.1s. The tracking results are shown in Table VIII.

As we can see, in this case, the robust filters (i.e., the IMM-RS, IMM-R, and DRIMM) are not preferable. In contrast, the adaptive filters (i.e., the IMM-B and IMM-M) outperform the nominal IMM-N filter because the uncertainty in the TPM is reduced. In particular, the IMM-C filter significantly works best. This is because the car is a slowly-maneuvering object and seldom maneuvers: the nearly-constant velocity (CV)

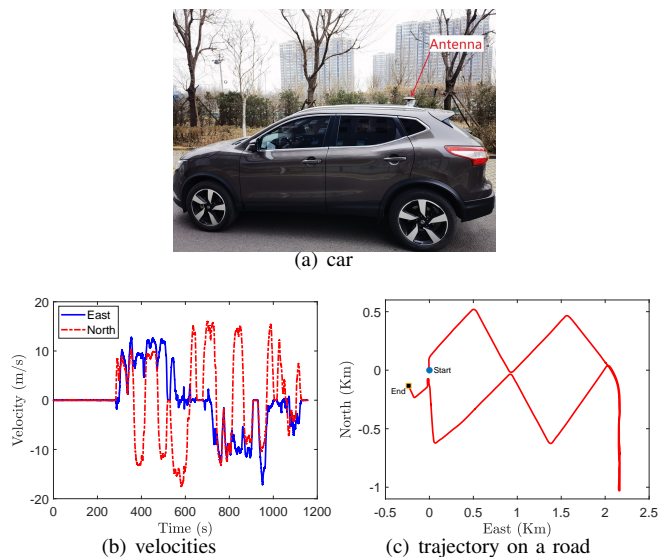


Fig. 1. Car, its real-time velocities, and its trajectory (data credit: UniStrong).

TABLE VIII  
TRACKING RESULTS OF THE CAR (SEE TABLE I FOR NOTES)

Filter	RMSE	Avg Time	Filter	RMSE	Avg Time
IMM-N	2.30	2.21e-05	IMM-B	2.26	2.86e-05
IMM-M	2.02	5.75e-05	IMM-C	<b>1.72</b>	2.28e-05
IMM-RS	2.30	8.22e-05	DRIMM	2.30	5.07e-05
IMM-R	2.29	2.34e-05			

model where  $\bar{a}_k = 0$  is the dominating model (N.B.: the trajectory consists of piece-wise straight-line segments). To clarify further, small uncertainties exist in the nominal model set and the CV model matches well most of the time.

2) *Track A Highly-Maneuvering Drone*: In this experiment, the data are collected by a highly-maneuvering quadrotor drone that carries a GPS solution and an RTK solution. The commercial model of the drone is Matrice-300-RTK produced by DJI Co., Ltd., Shenzhen, China; see <https://www.dji.com/>. The drone flies following round trajectories in the air over an open playground, and the flying speed is about 6m/s when collecting data. Parts of the drone's real-time velocities and its trajectory are shown in Fig. 2.

Without loss of generality, we study the tracking problem in the east axis (in the east-north-up coordinate). We suppose that there are three nominal values for acceleration in the east axis:  $\bar{a}_{j,k} = [0, 2.5, -2.5]$ . The diagonal elements of the nominal TPM are set to 0.8s and the non-diagonal ones are set to 0.1s. The tracking results are shown in Table IX.

TABLE IX  
TRACKING RESULTS OF THE DRONE (SEE TABLE I FOR NOTES)

Filter	RMSE	Avg Time	Filter	RMSE	Avg Time
IMM-N	0.99	4.85e-05	IMM-B	1.06	6.25e-05
IMM-M	1.00	1.00e-04	IMM-C	1.90	5.02e-05
IMM-RS	0.99	1.65e-04	DRIMM	<b>0.86</b>	1.20e-03
IMM-R	0.89	8.84e-05			

As we can see, in this case, the proposed DRIMM filter outperforms other filters. This is because the drone continu-

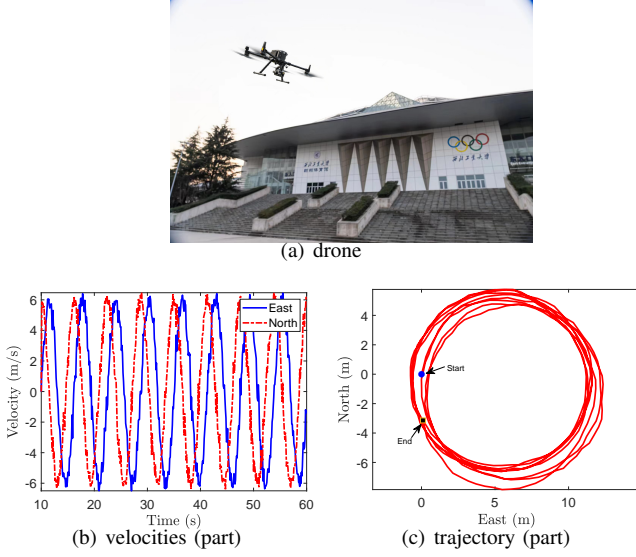


Fig. 2. Drone, its real-time velocities, and traveling trajectory, starting from 10s and ending at 60s (data credit: Competitive Robot Base, Northwestern Polytechnical University, Xi'An, China).

ously maneuvers all the time (i.e., the operating models are frequently switching), and therefore, large uncertainties exist in both the nominal model set and the TPM. Further, we improve the flying speed of the drone from 6m/s to 12m/s. (The maximum flying speed of the used drone is 23m/s; see <https://www.dji.com/>.) The corresponding tracking results are shown in Table X. The results suggest that, since the modeling uncertainties in both the nominal model set and the TPM are larger, the advantage of the proposed DRIMM filter becomes much more significant. In this case, the IMM-C filter (which performs best in the car-tracking case) even diverges because it is overly confident about the model that has the highest estimated model probability.

TABLE X  
TRACKING RESULTS OF THE DRONE (12M/S; SEE TABLE I FOR NOTES)

Filter	RMSE	Avg Time	Filter	RMSE	Avg Time
IMM-N	2.78	4.81e-05	IMM-B	2.68	6.12e-05
IMM-M	2.80	1.26e-04	IMM-C	Diverged	
IMM-RS	2.78	1.72e-04	DRIMM	<b>2.25</b>	9.22e-04
IMM-R	2.72	7.30e-05			

3) *Closing Remarks:* For detailed parameter settings, see shared source codes. The influences of the parameters used in filters are consistent with conclusions in simulated experiments in Section VI (e.g., sizes of ambiguity sets in Subsection VI-A5), and therefore, related discussions are omitted. One may use shared source codes to verify this claim.

## VII. CONCLUSIONS

This paper studies the robust state estimation problem for jump Markov linear systems subject to several modeling uncertainties such as 1) the uncertainty of the model transition probability matrix, 2) the uncertainties of nominal candidate models in the nominal model set, and/or 3) the incompleteness of the nominal model set. A unified distributionally robust

interactive multiple model (DRIMM) filter is proposed to withstand these model uncertainties. The following main points have to be highlighted.

- 1) When there exist uncertainties only in TPMs, the IMM filter with Bayesian estimation (i.e., the IMM-B method) and the IMM filter with maximum likelihood estimation (i.e., the IMM-M method) that aim to estimate the unknown true TPMs are recommended. All robust methods should be the last choices. However, the recommendation is conditioned on a sufficiently long time horizon (i.e., sufficient measurements are available so that the TPM can be estimated to a satisfactory accuracy level).<sup>6</sup>
- 2) When there exists a dominating operating model over time (i.e., the systems' operating modes do not frequently switch), the IMM filter with compensation (i.e., the IMM-C method) is recommended, for example, in slowly-maneuvering target tracking.
- 3) When there exist uncertainties in nominal candidate models or when the nominal model set is not complete, the proposed DRIMM method is recommended. This is because the DRIMM method is able to flexibly respond to all specified types of uncertainty. Also, it tends to be less conservative than the risk-sensitive IMM filter (i.e., the IMM-RS method).

Unfortunately, the performance of the DRIMM filter depends heavily on the size parameter(s) of the associated ambiguity set, and the convincing tuning method for the size parameter(s) is yet to be found. Therefore, future research following this work is expected to address this open issue.

## APPENDIX A

### OPTIMAL ESTIMATES AND ERROR COVARIANCES

We have

$$\hat{\mathbf{x}}_j = \bar{\mathbf{x}}_j + \mathbf{M}_j \mathbf{H}_j^\top (\mathbf{H}_j \mathbf{M}_j \mathbf{H}_j^\top + \mathbf{R}_j)^{-1} (\mathbf{y} - \mathbf{H}_j \bar{\mathbf{x}}_j), \quad (38)$$

$$\bar{\mathbf{P}}_j = \mathbf{M}_j - \mathbf{M}_j \mathbf{H}_j^\top (\mathbf{H}_j \mathbf{M}_j \mathbf{H}_j^\top + \mathbf{R}_j)^{-1} \mathbf{H}_j \mathbf{M}_j, \quad (39)$$

$$\hat{\mathbf{x}}_j = \mathbf{c}_{j,x} + \boldsymbol{\Sigma}_{j,x} \mathbf{H}_j^\top (\mathbf{H}_j \boldsymbol{\Sigma}_{j,x} \mathbf{H}_j^\top + \boldsymbol{\Sigma}_{j,v})^{-1} (\mathbf{y} - \mathbf{H}_j \mathbf{c}_{j,x} - \mathbf{c}_{j,v}), \quad (40)$$

and

$$\mathbf{P}_j = \boldsymbol{\Sigma}_{j,x} - \boldsymbol{\Sigma}_{j,x} \mathbf{H}_j^\top (\mathbf{H}_j \boldsymbol{\Sigma}_{j,x} \mathbf{H}_j^\top + \boldsymbol{\Sigma}_{j,v})^{-1} \mathbf{H}_j \boldsymbol{\Sigma}_{j,x}, \quad (41)$$

respectively. For detailed derivations of (38), (39), (40), and (41), refer to the standard Kalman filtering theory.

## APPENDIX B

### PROOF OF PROPOSITION 1

According to the strong min-max property in (14), (11) is equivalent to (13). Since for every  $\mathbb{P} \in \mathcal{F}_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta})$ , the associated optimal estimator  $\hat{\mathbf{x}} \in \mathcal{H}_{\mathbf{y}}$  is uniquely determined by  $\hat{\mathbf{x}} = \mathbb{E}(\mathbf{x}|\mathbf{y})$ , (13) reduces to finding the worst-case conditional distribution  $\mathbb{P}_{\mathbf{x}|\mathbf{y}=\mathbf{y}}^* \in \mathcal{F}_{\mathbf{x},\mathbf{y}=\mathbf{y}}(\boldsymbol{\theta})$  that solves

$$\max_{\mathbb{P} \in \mathcal{F}_{\mathbf{x},\mathbf{y}=\mathbf{y}}(\boldsymbol{\theta})} \text{Tr} \mathbb{E} \left\{ [\mathbf{x} - \phi(\mathbf{y})][\mathbf{x} - \phi(\mathbf{y})]^\top \middle| \mathbf{y} = \mathbf{y} \right\}, \quad (42)$$

<sup>6</sup>One may use shared source codes to verify this claim.

where the ambiguity set  $\mathcal{F}_{\mathbf{x}, \mathbf{y}=\mathbf{y}}(\theta)$  contains all conditional distributions of  $\mathbf{x}$  given  $\mathbf{y} = \mathbf{y}$ . Since every candidate distribution  $\mathbb{P}_{\mathbf{x}|\mathbf{y}=\mathbf{y}}$  in  $\mathcal{F}_{\mathbf{x}, \mathbf{y}=\mathbf{y}}(\theta)$  is parameterized by  $\{\omega_j, \mathbf{c}_{j,\mathbf{x}}, \Sigma_{j,\mathbf{x}}, \mathbf{c}_{j,\mathbf{y}}, \Sigma_{j,\mathbf{y}}\}_{j \in [N]}$ , by recalling (3), (40), (41), (17), and (18), problem (42) can be transformed into (19).  $\square$

#### APPENDIX C

##### USE WASSERSTEIN DISTANCE BETWEEN $\mu$ AND $\bar{\mu}$

Below we discuss the case where the statistical similarity measure between  $\mu$  and  $\bar{\mu}$  is defined by the Wasserstein metric. The distributionally robust Bayesian problem (25) is particularized to

$$\begin{aligned} \max_{\mu} \quad & -\mu^\top \mathbf{A}\mu + \mathbf{b}^\top \mu \\ \text{s.t.} \quad & \begin{cases} \min_{\mathbf{Q}} \sum_{i=1}^N \sum_{j=1}^N C_{ij} Q_{ij} \leq \theta_0, \\ \sum_{j=1}^N Q_{ij} = \bar{\mu}_i, \quad \forall i \in [N], \\ \sum_{i=1}^N Q_{ij} = \mu_j, \quad \forall j \in [N], \\ Q_{ij} \geq 0, \quad \forall i, j \in [N], \end{cases} \end{aligned} \quad (43)$$

where  $C_{ij}$  measures the difference between the two distributions  $\mathbb{P}_{i,\mathbf{x},\mathbf{y}}$  and  $\mathbb{P}_{j,\mathbf{x},\mathbf{y}}$  (which can be defined by the Wasserstein distance, the KL divergence, or whatever suitable), and the matrix  $\mathbf{Q} := \{Q_{ij}\}_{i,j \in [N]}$  can be regarded as a discrete joint distribution whose marginals are  $\mu$  and  $\bar{\mu}$ . Note that  $\sum_{j=1}^N \mu_j = 1 = \sum_{j=1}^N \sum_{i=1}^N Q_{ij}$  is implicitly admitted because  $\sum_{i=1}^N \sum_{j=1}^N Q_{ij} = \sum_{i=1}^N \bar{\mu}_i = 1$ . Note also that the minimization operator in the first constraint can be dropped because it is redundant for this problem. For two Gaussian distributions  $\mathcal{N}_d(\mathbf{c}_i, \Sigma_i)$  and  $\mathcal{N}_d(\mathbf{c}_j, \Sigma_j)$ , the Wasserstein distance between them

is  $C_{ij}^W = \sqrt{\|\mathbf{c}_i - \mathbf{c}_j\|^2 + \text{Tr} \left[ \Sigma_i + \Sigma_j - 2 \left( \Sigma_i^{\frac{1}{2}} \Sigma_j \Sigma_i^{\frac{1}{2}} \right)^{\frac{1}{2}} \right]}$  and the KL divergence between them is  $C_{ij}^{\text{KL}} = \frac{1}{2} \left[ \|\mathbf{c}_i - \mathbf{c}_j\|_{\Sigma_j^{-1}}^2 + \text{Tr} [\Sigma_j^{-1} \Sigma_i - \mathbf{I}] + \ln \det (\Sigma_i^{-1} \Sigma_j) \right]$ .

By plugging  $\sum_{i=1}^N Q_{ij} = \mu_j, \forall j \in [N]$  into the objective function, the problem (43) can be explicitly written as

$$\begin{aligned} \max_{\mathbf{Q}} \quad & -\sum_{i=1}^N \sum_{j=1}^N \sum_{m=1}^N \sum_{n=1}^N Q_{mi} A_{ij} Q_{nj} + \sum_{i=1}^N \sum_{j=1}^N b_j Q_{ij} \\ \text{s.t.} \quad & \begin{cases} \sum_{i=1}^N \sum_{j=1}^N C_{ij} Q_{ij} \leq \theta_0, \\ \sum_{j=1}^N Q_{ij} = \bar{\mu}_i, \quad \forall i \in [N], \\ Q_{ij} \geq 0, \quad \forall i, j \in [N]. \end{cases} \end{aligned} \quad (44)$$

The problem (44) is quadratic in the matrix-valued variable  $\mathbf{Q}$ . Therefore, solving (44) is difficult because existing general-purpose commercial solvers cannot be used. As a result, a specifically efficient solution method for (44) or (43) is expected. We propose to use the Frank-Wolfe method [41]. The Frank-Wolfe method is an iterative method. At each iteration, it linearizes the nonlinear objective of an optimization problem and uses the objective-linearized sub-problem to find a feasible direction along which the current solution can be improved. Specifically, the Frank-Wolfe method can construct a sequence  $\{\mu^{(r)}\}_{r=0,1,2,\dots}$  such that it converges to the optimal solution of (43) as  $r \rightarrow \infty$  where  $r$  denotes the iteration count.

*Proposition 5:* Let  $r = 1, 2, \dots$  be the iteration count. Let  $\mathbf{c}^{(r)} := -2\mathbf{A}\mu^{(r)} + \mathbf{b}$  denote the gradient of the objective of (43) at  $\mu^{(r)}$ . Construct the following iteration process

$$\mu^{(r+1)} = \mu^{(r)} + \beta_r \cdot (\mathbf{s}^{(r)} - \mu^{(r)}) \quad (45)$$

where  $\mathbf{s}^{(r)} := [s_1^{(r)}, s_2^{(r)}, \dots, s_N^{(r)}]^\top$ ,  $s_j^{(r)} := \sum_{i=1}^N Q_{ij}^{*(r)}$ ,  $\forall j \in [N]$ ,

$$\begin{aligned} \mathbf{Q}^{*(r)} &:= \arg\max_{\mathbf{Q}} \sum_{i=1}^N \sum_{j=1}^N \mathbf{c}_j^{(r)} \cdot Q_{ij} \\ \text{s.t.} \quad & \begin{cases} \sum_{i=1}^N \sum_{j=1}^N C_{ij} Q_{ij} \leq \theta_0, \\ \sum_{j=1}^N Q_{ij} = \bar{\mu}_i, \quad \forall i \in [N], \\ Q_{ij} \geq 0, \quad \forall i, j \in [N], \end{cases} \end{aligned} \quad (46)$$

and the step size  $\beta_r := \frac{2}{r+1}, \forall r$ . Then  $\mu^{(r)}$  converges to an optimal solution of (43), as  $r \rightarrow \infty$ . [Note that the solution of (43) might not be unique.]

*Proof:* See Appendix D.  $\square$

Note that the sub-problem in (46) is a  $N^2$ -variate linear program that can be solved by the simplex method. However, the sub-problem (46) still introduces higher computational burdens to the solution method under the Wasserstein distance compared to the solution method under the KL divergence. To be specific, the computational complexity of the simplex method for (46) in the worst case is exponential in  $N^2$ ; no better theoretical results are reported for the generic simplex method although it performs well empirically. However, the root-finding sub-problem in Algorithm 2 in the online supplementary materials is just polynomial in  $N$  if Newton's method is employed.

*Remark 5:* Note that in the Wasserstein case, (26) in Line 8 of Algorithm 1 should be replaced with (43).  $\square$

*Remark 6:* Different statistical similarity measures define different ambiguity sets (e.g., different geometric shapes), and different ambiguity sets further lead to different robust state estimation results [30]. Therefore, in practice, if computational powers allow, one can try all possible ambiguity sets to obtain better performance on specific problems.  $\square$

#### APPENDIX D

##### PROOF OF PROPOSITION 5

This proposition specifies the Frank-Wolfe method [41] for the problem (43). Let  $\mathbf{s}^{(r)}$  solve the objective-linearized sub-problem of (43)

$$\begin{aligned} \mathbf{s}^{(r)} &= \arg\max_{\mu} [\mathbf{c}^{(r)}]^\top \mu \\ \text{s.t.} \quad & \begin{cases} \sum_{i=1}^N \sum_{j=1}^N C_{ij} Q_{ij} \leq \theta_0, \\ \sum_{j=1}^N Q_{ij} = \bar{\mu}_i, \quad \forall i \in [N], \\ \sum_{i=1}^N Q_{ij} = \mu_j, \quad \forall j \in [N], \\ Q_{ij} \geq 0, \quad \forall i, j \in [N]. \end{cases} \end{aligned} \quad (47)$$

By plugging in  $\mu_j := \sum_{i=1}^N Q_{ij}, \forall j \in [N]$  into the objective, (47) can be written as (46), which is a linear program that can be efficiently solved by the simplex method. In (45), the vector  $\mathbf{s}^{(r)} - \mu^{(r)}$  gives a feasible direction along which the current solution  $\mu^{(r)}$  can be improved to  $\mu^{(r+1)}$ . The step size  $\beta_r$  can be  $\beta_r := \frac{2}{r+1}, \forall r$  [41]. Since the objective function

of (43) is concave and continuously differentiable, and the feasible region of (43) is a simplex (thus compact and convex) after dropping the minimization operator in the first constraint, according to [41, Theorem 1], this proposition holds.  $\square$

#### APPENDIX E PROOF OF THEOREM 1

We first prove the first part: for every given  $\Sigma_{j,v}$ , if  $\Sigma_1 \succeq \Sigma_2$ , we have  $\text{Tr } P_j(\Sigma_1, \Sigma_{j,v}) \geq \text{Tr } P_j(\Sigma_2, \Sigma_{j,v})$ . We consider two optimization problems:

$$U_1 = \underset{U}{\text{argmax}} \text{Tr } U$$

$$\text{s.t.} \begin{cases} \begin{bmatrix} \Sigma_1 - U & \Sigma_1 H_j^\top \\ H_j \Sigma_1 & H_j \Sigma_1 H_j^\top + \Sigma_{j,v} \end{bmatrix} \succeq 0, \\ U \succeq 0. \end{cases} \quad (48)$$

$$U_2 = \underset{U}{\text{argmax}} \text{Tr } U$$

$$\text{s.t.} \begin{cases} \begin{bmatrix} \Sigma_2 - U & \Sigma_2 H_j^\top \\ H_j \Sigma_2 & H_j \Sigma_2 H_j^\top + \Sigma_{j,v} \end{bmatrix} \succeq 0, \\ U \succeq 0. \end{cases} \quad (49)$$

Since  $H_j \Sigma_1 H_j^\top + \Sigma_{j,v} \succ 0$  and  $H_j \Sigma_2 H_j^\top + \Sigma_{j,v} \succ 0$ , by Schur complement, we have  $U_1 = \Sigma_1 - \Sigma_1 H_j^\top (H_j \Sigma_1 H_j^\top + \Sigma_{j,v})^{-1} H_j \Sigma_1$  and  $U_2 = \Sigma_2 - \Sigma_2 H_j^\top (H_j \Sigma_2 H_j^\top + \Sigma_{j,v})^{-1} H_j \Sigma_2$ . On the other hand,

$$\begin{bmatrix} U & 0 \\ 0 & 0 \end{bmatrix} \preceq \begin{bmatrix} \Sigma_2 & \Sigma_2 H_j^\top \\ H_j \Sigma_2 & H_j \Sigma_2 H_j^\top + \Sigma_{j,v} \end{bmatrix}$$

$$\preceq \begin{bmatrix} \Sigma_1 & \Sigma_1 H_j^\top \\ H_j \Sigma_1 & H_j \Sigma_1 H_j^\top + \Sigma_{j,v} \end{bmatrix},$$

which implies that the feasible region of (48) is larger than that of (49). Therefore,  $U_1 \succeq U_2$ , and  $P_j(\Sigma_1, \Sigma_{j,v}) - P_j(\Sigma_2, \Sigma_{j,v}) = U_1 - U_2 \succeq 0$ .

Next, we prove the second part: for every given  $\Sigma_{j,x}$ , if  $\Sigma_1 \succeq \Sigma_2$ , we have  $\text{Tr } P_j(\Sigma_{j,x}, \Sigma_1) \geq \text{Tr } P_j(\Sigma_{j,x}, \Sigma_2)$ . Since  $P_j(\Sigma_{j,x}, \Sigma_1) = \Sigma_{j,x} - \Sigma_{j,x} H_j^\top (H_j \Sigma_{j,x} H_j^\top + \Sigma_1)^{-1} H_j \Sigma_{j,x}$ ,  $P_j(\Sigma_{j,x}, \Sigma_2) = \Sigma_{j,x} - \Sigma_{j,x} H_j^\top (H_j \Sigma_{j,x} H_j^\top + \Sigma_2)^{-1} H_j \Sigma_{j,x}$ , the conclusion is obvious due to basic algebra.  $\square$

#### APPENDIX F

##### QUANTIFY UNCERTAINTY USING WASSERSTEIN SETS

If we use the Wasserstein ambiguity sets for  $\mathbb{P}_{j,x,y}$  as in [13], the problem (33) is particularized into

$$\max_{\Sigma_{j,x}, \Sigma_{j,v}} \text{Tr}[\Sigma_{j,x} + \Sigma_{j,v}]$$

$$\text{s.t.} \begin{cases} \Sigma_{j,x} H_j^\top = K_j \cdot (H_j \Sigma_{j,x} H_j^\top + \Sigma_{j,v}), \\ \sqrt{\text{Tr}[\Sigma_{j,x} + M_j - 2(M_j^{\frac{1}{2}} \Sigma_{j,x} M_j^{\frac{1}{2}})^{\frac{1}{2}}]} \leq \theta_{j,x}, \\ \sqrt{\text{Tr}[\Sigma_{j,v} + R_j - 2(R_j^{\frac{1}{2}} \Sigma_{j,v} R_j^{\frac{1}{2}})^{\frac{1}{2}}]} \leq \theta_{j,v}, \\ \Sigma_{j,x} \succeq 0, \Sigma_{j,v} \succeq 0. \end{cases} \quad (50)$$

where  $\theta_{j,x} \geq 0$  and  $\theta_{j,v} \geq 0$ . Problem (50) is a nonlinear positive semi-definite program (SDP). However, it has a linear reformulation, which can be efficiently solved.

**Proposition 6:** The nonlinear positive SDP (50) can be reformulated into a linear positive SDP

$$\max_{\Sigma_{j,x}, \Sigma_{j,v}, V_{j,x}, V_{j,v}} \text{Tr}[\Sigma_{j,x} + \Sigma_{j,v}]$$

$$\text{s.t.} \begin{cases} \Sigma_{j,x} H_j^\top = K_j \cdot (H_j \Sigma_{j,x} H_j^\top + \Sigma_{j,v}), \\ \text{Tr}[\Sigma_{j,x} + M_j - 2V_{j,x}] \leq \theta_{j,x}^2, \\ \begin{bmatrix} M_j^{\frac{1}{2}} \Sigma_{j,x} M_j^{\frac{1}{2}} & V_{j,x} \\ V_{j,x} & I \end{bmatrix} \succeq 0, \\ \text{Tr}[\Sigma_{j,v} + R_j - 2V_{j,v}] \leq \theta_{j,v}^2, \\ \begin{bmatrix} R_j^{\frac{1}{2}} \Sigma_{j,v} R_j^{\frac{1}{2}} & V_{j,v} \\ V_{j,v} & I \end{bmatrix} \succeq 0, \\ \Sigma_{j,x} \succeq 0, \Sigma_{j,v} \succeq 0, V_{j,x} \succeq 0, V_{j,v} \succeq 0. \end{cases} \quad (51)$$

*Proof:* By introducing  $M_j^{\frac{1}{2}} \Sigma_{j,x} M_j^{\frac{1}{2}} \succeq V_{j,x}^2$  where  $V_{j,x} \succeq 0$  and using Schur complement, the nonlinear constraint  $\sqrt{\text{Tr}[\Sigma_{j,x} + M_j - 2(M_j^{\frac{1}{2}} \Sigma_{j,x} M_j^{\frac{1}{2}})^{\frac{1}{2}}]} \leq \theta_{j,x}$  in (50) can be reformulated into a linear equivalent

$$\begin{cases} \text{Tr}[\Sigma_{j,x} + M_j - 2V_{j,x}] \leq \theta_{j,x}^2, \\ \begin{bmatrix} M_j^{\frac{1}{2}} \Sigma_{j,x} M_j^{\frac{1}{2}} & V_{j,x} \\ V_{j,x} & I \end{bmatrix} \succeq 0, \\ V_{j,x} \succeq 0. \end{cases}$$

The constraint  $\sqrt{\text{Tr}[\Sigma_{j,v} + R_j - 2(R_j^{\frac{1}{2}} \Sigma_{j,v} R_j^{\frac{1}{2}})^{\frac{1}{2}}]} \leq \theta_{j,v}$  in (50) can be handled similarly. This completes the proof.  $\square$

**Remark 7:** Note that in the Wasserstein case, (34) in Line 5 of Algorithm 1 should be replaced with (51).  $\square$

#### APPENDIX G PROOF OF THEOREM 2

To begin with, we drop the first constraint  $\Sigma_{j,x} H_j^\top = K_j \cdot (H_j \Sigma_{j,x} H_j^\top + \Sigma_{j,v})$  and solve the following relaxed problem:

$$\max_{\Sigma_{j,x}, \Sigma_{j,v}} \text{Tr}[\Sigma_{j,x} + \Sigma_{j,v}]$$

$$\text{s.t.} \begin{cases} (1 - \theta_{j,x}) M_j \preceq \Sigma_{j,x} \preceq (1 + \theta_{j,x}) M_j, \\ (1 - \theta_{j,v}) R_j \preceq \Sigma_{j,v} \preceq (1 + \theta_{j,v}) R_j, \\ \Sigma_{j,x} \succeq 0, \Sigma_{j,v} \succeq 0. \end{cases}$$

The solution of this relaxed problem is obviously given in (35). By verifying that the solution (35) also satisfies the constraint  $\Sigma_{j,x} H_j^\top = K_j \cdot (H_j \Sigma_{j,x} H_j^\top + \Sigma_{j,v})$  where  $K_j$  is defined in (29), we complete the proof.  $\square$

#### REFERENCES

- [1] X. R. Li and V. P. Jilkov, "Survey of maneuvering target tracking. part v. multiple-model methods," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 41, no. 4, pp. 1255–1321, 2005.
- [2] —, "Survey of maneuvering target tracking. part i. dynamic models," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 39, no. 4, pp. 1333–1364, 2003.
- [3] Y. Bar-Shalom, S. Challa, and H. A. Blom, "IMM estimator versus optimal estimator for hybrid systems," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 41, no. 3, pp. 986–991, 2005.
- [4] Y. Bar-Shalom, X. R. Li, and T. Kirubarajan, *Estimation with Applications to Tracking and Navigation*. John Wiley & Sons, 2001.
- [5] H. A. Blom and Y. Bar-Shalom, "The interacting multiple model algorithm for systems with Markovian switching coefficients," *IEEE Transactions on Automatic Control*, vol. 33, no. 8, pp. 780–783, 1988.



- [6] E. Mazor, A. Averbuch, Y. Bar-Shalom, and J. Dayan, "Interacting multiple model methods in target tracking: a survey," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 34, no. 1, pp. 103–123, 1998.
- [7] V. P. Jilkov and X. R. Li, "Online Bayesian estimation of transition probabilities for markovian jump systems," *IEEE Transactions on Signal Processing*, vol. 52, no. 6, pp. 1620–1630, 2004.
- [8] U. Orguner and M. Demirekler, "Maximum likelihood estimation of transition probabilities of jump Markov linear systems," *IEEE Transactions on Signal Processing*, vol. 56, no. 10, pp. 5093–5108, 2008.
- [9] S. Zhao and F. Liu, "Recursive estimation for Markov jump linear systems with unknown transition probabilities: A compensation approach," *Journal of the Franklin Institute*, vol. 353, no. 7, pp. 1494–1517, 2016.
- [10] Y. Ma, S. Zhao, and B. Huang, "Multiple-model state estimation based on variational Bayesian inference," *IEEE Transactions on Automatic Control*, vol. 64, no. 4, pp. 1679–1685, 2018.
- [11] A. Doucet and B. Ristic, "Recursive state estimation for multiple switching models with unknown transition probabilities," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 38, no. 3, pp. 1098–1104, 2002.
- [12] U. Orguner and M. Demirekler, "An online sequential algorithm for the estimation of transition probabilities for jump Markov linear systems," *Automatica*, vol. 42, no. 10, pp. 1735–1744, 2006.
- [13] S. Wang, Z. Wu, and A. Lim, "Robust state estimation for linear systems under distributional uncertainty," *IEEE Transactions on Signal Processing*, 2021.
- [14] S. Wang and Z.-S. Ye, "Distributionally robust state estimation for linear systems subject to uncertainty and outlier," *IEEE Transactions on Signal Processing*, vol. 70, pp. 452–467, 2021.
- [15] F. Caron, M. Davy, A. Doucet, E. Duflos, and P. Vanheeghe, "Bayesian inference for linear dynamic models with dirichlet process mixtures," *IEEE Transactions on Signal Processing*, vol. 56, no. 1, pp. 71–84, 2007.
- [16] C. Magnan, A. Giremus, E. Grivel, L. Ratton, and B. Joseph, "Bayesian non-parametric methods for dynamic state-noise covariance matrix estimation: Application to target tracking," *Signal Processing*, vol. 127, pp. 135–150, 2016.
- [17] R. Guo, M. Shen, D. Huang, X. Yin, and L. Xu, "Recursive estimation of transition probabilities for jump Markov linear systems under minimum Kullback–Leibler divergence criterion," *IET Control Theory & Applications*, vol. 9, no. 17, pp. 2491–2499, 2015.
- [18] L. Zhang and E.-K. Boukas, "Mode-dependent  $\mathcal{H}_\infty$  filtering for discrete-time Markovian jump linear systems with partly unknown transition probabilities," *Automatica*, vol. 45, no. 6, pp. 1462–1467, 2009.
- [19] X. Li, J. Lam, H. Gao, and J. Xiong, " $\mathcal{H}_\infty$  and  $\mathcal{H}_2$  filtering for linear systems with uncertain Markov transitions," *Automatica*, vol. 67, pp. 252–266, 2016.
- [20] C. E. de Souza, A. Trofino, and K. A. Barbosa, "Mode-independent  $\mathcal{H}_\infty$  filters for Markovian jump linear systems," *IEEE Transactions on Automatic Control*, vol. 51, no. 11, pp. 1837–1841, 2006.
- [21] U. Orguner and M. Demirekler, "Risk-sensitive filtering for jump markov linear systems," *Automatica*, vol. 44, no. 1, pp. 109–118, 2008.
- [22] S. Zhao and F. Liu, "An adaptive risk-sensitive filtering method for markov jump linear systems with uncertain parameters," *Journal of the Franklin Institute*, vol. 349, no. 6, pp. 2047–2064, 2012.
- [23] M. J. Schervish, *Theory of Statistics*. Springer Science & Business Media, 2012.
- [24] B. D. Anderson and J. B. Moore, *Optimal Filtering*. Prentice-Hall, 1979.
- [25] T. Kailath, A. H. Sayed, and B. Hassibi, *Linear Estimation*. Prentice Hall, 2000.
- [26] D. Simon, *Optimal State Estimation: Kalman,  $H_\infty$ , and Nonlinear Approaches*. John Wiley & Sons, 2006.
- [27] B. Hassibi, A. H. Sayed, and T. Kailath, "Linear estimation in Krein spaces. I. theory," *IEEE Transactions on Automatic Control*, vol. 41, no. 1, pp. 18–33, 1996.
- [28] B. C. Levy and R. Nikoukhan, "Robust state space filtering under incremental model perturbations subject to a relative entropy tolerance," *IEEE Transactions on Automatic Control*, vol. 58, no. 3, pp. 682–695, 2013.
- [29] S. Shafieezadeh Abadeh, V. A. Nguyen, D. Kuhn, and P. M. Mohajerin Esfahani, "Wasserstein distributionally robust kalman filtering," *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [30] S. Wang, "Distributionally robust state estimation," Ph.D. dissertation, National University of Singapore, 1 2022.
- [31] A. Ben-Tal, D. den Hertog, A. D. Waegenaere, B. Melenberg, and G. Rennen, "Robust solutions of optimization problems affected by uncertain probabilities," *Management Science*, vol. 59, pp. 341–357, 2 2013.
- [32] P. Mohajerin Esfahani and D. Kuhn, "Data-driven distributionally robust optimization using the wasserstein metric: Performance guarantees and tractable reformulations," *Mathematical Programming*, vol. 171, no. 1–2, pp. 115–166, 2018.
- [33] D. Kuhn, P. M. Esfahani, V. A. Nguyen, and S. Shafieezadeh-Abadeh, "Wasserstein distributionally robust optimization: Theory and applications in machine learning," in *Operations Research & Management Science in the Age of Analytics*. INFORMS, 2019, pp. 130–166.
- [34] R. Gao and A. Kleywegt, "Distributionally robust stochastic optimization with wasserstein distance," *Mathematics of Operations Research*, vol. 48, no. 2, pp. 603–655, 2022.
- [35] R. Zhu, H. Wei, and X. Bai, "Wasserstein metric based distributionally robust approximate framework for unit commitment," *IEEE Transactions on Power Systems*, vol. 34, no. 4, pp. 2991–3001, 2019.
- [36] B. Zhou, G. Chen, T. Huang, Q. Song, and Y. Yuan, "Planning pev fast-charging stations using data-driven distributionally robust optimization approach based on  $\phi$ -divergence," *IEEE Transactions on Transportation Electrification*, vol. 6, no. 1, pp. 170–180, 2020.
- [37] S. Boyd, S. P. Boyd, and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.
- [38] G. Wang, "ML estimation of transition probabilities in jump Markov systems via convex optimization," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 46, no. 3, pp. 1492–1502, 2010.
- [39] S. Wang, "Distributionally robust state estimation for nonlinear systems," *IEEE Transactions on Signal Processing*, vol. 70, pp. 4408–4423, 2022.
- [40] P. Teunissen and A. Khodabandeh, "Review and principles of PPP-RTK methods," *Journal of Geodesy*, vol. 89, no. 3, pp. 217–240, 2015.
- [41] M. Jaggi, "Revisiting Frank-Wolfe: Projection-free sparse convex optimization," in *International Conference on Machine Learning*. PMLR, 2013, pp. 427–435.



**Dr. Shixiong Wang** received the B.Eng. degree in detection, guidance, and control technology, and the M.Eng. degree in systems and control engineering from the School of Electronics and Information, Northwestern Polytechnical University, China, in 2016 and 2018, respectively. He received his Ph.D. degree from the Department of Industrial Systems Engineering and Management, National University of Singapore, Singapore, in 2022.

He is currently a Postdoctoral Research Associate with the Intelligent Transmission and Processing Laboratory, Imperial College London, London, United Kingdom, from May 2023. He was a Postdoctoral Research Fellow with the Institute of Data Science, National University of Singapore, Singapore, from March 2022 to March 2023.

His research interest includes statistics and optimization theories with applications in signal processing (especially optimal estimation theory), machine learning (especially generalization error theory), and control technology.



## Supplementary Materials

### APPENDIX H ROBUSTIFICATION OVER PRIOR MODEL PROBABILITY

In this appendix, we discuss the case where we solve the genuine problem (24) with respect to  $\omega$ , rather than the simplified case with respect to  $\mu$ . In consideration of the high computational complexity, the genuine problem (24) with respect to  $\omega$  is not investigated in the main body of the paper.

*Proposition 7:* If the model set is exact and only the prior model probability vector  $\omega$  is uncertain [i.e., the special ambiguity set (23) is investigated], the reformulated distributionally robust Bayesian estimation problem (24) can be further reformulated into a tractable quadratic fractional program

$$\begin{aligned} \max_{\omega} \quad & -\frac{\omega^\top (CAC - pb^\top C)\omega}{\omega^\top pp^\top \omega} \\ \text{s.t.} \quad & \begin{cases} \sum_{j=1}^N \omega_j = 1, \\ \omega_j \geq 0, \\ \Delta_0(\omega, \bar{\omega}) \leq \theta_0, \end{cases} \quad \forall j \in [N], \end{aligned} \quad (52)$$

where  $p := [p_1(\mathbf{y}), p_2(\mathbf{y}), \dots, p_N(\mathbf{y})]^\top$  denotes the likelihoods of the candidate models given the measurement  $\mathbf{y}$  and  $C := \text{diag}(p)$  is a diagonal matrix whose diagonal entries are elements of  $p$ .

*Proof:* From (3), for every  $j \in [N]$ , we have  $\mu_j = \frac{\omega_j p_j(\mathbf{y})}{\sum_{j=1}^N \omega_j p_j(\mathbf{y})} = \frac{\omega_j p_j(\mathbf{y})}{\omega^\top p}$ , i.e.,  $\mu = [\mu_1, \mu_2, \dots, \mu_N]^\top = \frac{C\omega}{\omega^\top p}$ . Therefore, the problem (24) can be explicitly written as

$$\begin{aligned} \max_{\omega} \quad & -\left(\frac{C\omega}{\omega^\top p}\right)^\top A \left(\frac{C\omega}{\omega^\top p}\right) + b^\top \left(\frac{C\omega}{\omega^\top p}\right) \\ \text{s.t.} \quad & \begin{cases} \sum_{j=1}^N \omega_j = 1, \\ \omega_j \geq 0, \\ \Delta_0(\omega, \bar{\omega}) \leq \theta_0, \end{cases} \quad \forall j \in [N], \end{aligned} \quad (53)$$

which can be rearranged into the quadratic fractional program (52).  $\square$

The problem (52) can be written in a compact form

$$\max_{\omega \in \Omega} \frac{f_1(\omega)}{f_2(\omega)}, \quad (54)$$

where  $f_1(\omega) := -\omega^\top (CAC - pb^\top C)\omega$  denotes the numerator of the objective of (52),  $f_2(\omega) := \omega^\top pp^\top \omega$  the denominator of the objective of (52), and  $\Omega$  the feasible region of (52). One may verify that although  $f_2(\omega)$  is convex,  $f_1(\omega)$  is neither concave nor convex. However,  $f_1(\omega) \geq 0$  can be guaranteed because the objective of (19) is non-negative, as are those of (24) and (52). Complete (approximated) solutions to the problem (54) can be found in, e.g., [S1],<sup>7</sup> [S2],<sup>8</sup> where involved indefinite quadratic programs can be solved by the method in, e.g., [S3].<sup>9</sup> Numerically solving (54) is time-consuming due to the indefiniteness of  $f_1(\omega)$ . Therefore, in this paper, we do not proceed further for (54). Instead, we find a simplified alternative to the original problem (24) with respect to  $\mu$ . Interested readers may implement solution methods in, e.g., [S3], to solve (54) themselves.

### APPENDIX I SOLUTION TO (26)

The Lagrangian of (26) is

$$\begin{aligned} \min_{\lambda_0 \geq 0, \lambda_1} \max_{\mu} \quad & -\mu^\top A\mu + b^\top \mu + \lambda_1 \cdot (1 - \mathbf{1}^\top \mu) + \\ & \lambda_0 \cdot (\theta_0 - \mu^\top \ln \mu + \mu^\top \ln \bar{\mu}). \end{aligned} \quad (55)$$

For every  $\lambda_0 \geq 0$  and  $\lambda_1$ , the maximum  $\mu$  satisfies the first-order optimality condition:

$$-2A\mu + b - \lambda_1 \cdot \mathbf{1} + \lambda_0 \cdot (-\ln \mu - \mathbf{1} + \ln \bar{\mu}) = \mathbf{0}, \quad (56)$$

<sup>7</sup>[S1] W. Dinkelbach, "On nonlinear fractional programming," *Management Science*, vol. 13, no. 7, pp. 492–498, 1967.

<sup>8</sup>[S2] A. T. Phillips, *Quadratic Fractional Programming: Dinkelbach Method*. Boston, MA: Springer US, 2001, pp. 2107–2110. [Online]. Available: [https://doi.org/10.1007/0-306-48332-7\\_406](https://doi.org/10.1007/0-306-48332-7_406).

<sup>9</sup>[S3] A. Phillips and J. Rosen, "Guaranteed  $\epsilon$ -approximate solution for indefinite quadratic global minimization," *Naval Research Logistics (NRL)*, vol. 37, no. 4, pp. 499–514, 1990.

which transforms (55) to

$$\min_{\lambda_0 \geq 0, \lambda_1} \lambda_0 \theta_0 + \lambda_1 + \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu} + \lambda_0 \mathbf{1}^\top \boldsymbol{\mu}. \quad (57)$$

Since (26) is a convex program and  $\bar{\boldsymbol{\mu}}$  is a relative interior point in the feasible set, there does not exist duality gap between (26) and (57). Since (57) is convex, any first-order gradient-based method, e.g., projected gradient descent, is applicable to solve it. Let the objective of (57) be denoted as  $f(\boldsymbol{\lambda})$ . From (56), we have  $-2\mathbf{A} \frac{d\boldsymbol{\mu}}{d\lambda_0} = \ln \boldsymbol{\mu} + \mathbf{1} - \ln \bar{\boldsymbol{\mu}} + \lambda_0 \frac{1}{\boldsymbol{\mu}} \odot \frac{d\boldsymbol{\mu}}{d\lambda_0}$ , and  $-2\mathbf{A} \frac{d\boldsymbol{\mu}}{d\lambda_1} = \mathbf{1} + \lambda_0 \frac{1}{\boldsymbol{\mu}} \odot \frac{d\boldsymbol{\mu}}{d\lambda_1}$ , where  $\frac{1}{\boldsymbol{\mu}}$  means element-wise fraction, and  $\odot$  denotes the Hadamard product (i.e., the element-wise product). The gradient of the objective of (57) with respect to  $\lambda_0$  and  $\lambda_1$  are given by

$$\begin{aligned} \frac{\partial f(\boldsymbol{\lambda})}{\partial \lambda_0} &= \theta_0 + 2\boldsymbol{\mu}^\top \mathbf{A} \frac{d\boldsymbol{\mu}}{d\lambda_0} + \mathbf{1}^\top \boldsymbol{\mu} + \lambda_0 \mathbf{1}^\top \frac{d\boldsymbol{\mu}}{d\lambda_0} \\ &= \theta_0 - \boldsymbol{\mu}^\top \ln \boldsymbol{\mu} + \boldsymbol{\mu}^\top \ln \bar{\boldsymbol{\mu}}, \end{aligned} \quad (58)$$

and

$$\frac{\partial f(\boldsymbol{\lambda})}{\partial \lambda_1} = 1 + 2\boldsymbol{\mu}^\top \mathbf{A} \frac{d\boldsymbol{\mu}}{d\lambda_0} + \lambda_0 \mathbf{1}^\top \frac{d\boldsymbol{\mu}}{d\lambda_1} = 1 - \mathbf{1}^\top \boldsymbol{\mu}. \quad (59)$$

respectively. Hence, when the optimality of (57) reaches, i.e., when the gradients with respect to  $\lambda_0$  and  $\lambda_1$  vanish, we have  $1 = \sum_{j=1}^N \mu_j$  and  $\theta_0 = \sum_{j=1}^N \mu_j \cdot \ln \frac{\mu_j}{\bar{\mu}_j}$ . Specifically, it means  $\boldsymbol{\mu}$  is indeed a distribution summed to unit and all the robustness budget  $\theta_0$  has been used. In summary, the solution to (26) is summarized in Algorithm 2. Since (26) is a convex program, every iteration improves the objective.

---

**Algorithm 2** Solution to (26)

---

**Definition:**  $S$  as maximum allowed iteration steps and  $s$  the current iteration step;  $\alpha$  as step size;  $\epsilon$  as numerical precision threshold;  $\text{abs}(\cdot)$  returns absolute value.

**Remark:** Since (57) is convex, any initial values for  $\lambda_0 \geq 0$  and  $\lambda_1$  are acceptable. If early stopping is applied (i.e.,  $S$  is not sufficiently large for time-saving purpose), a normalization procedure is necessary to guarantee  $1 = \sum_j \mu_j$ .

**Input:**  $S, \alpha, \epsilon, \lambda_0, \lambda_1$

```

1:  $s \leftarrow 0$ ;
2: while true do
3:   // Update  $\boldsymbol{\mu}$ 
4:   Solve  $N$ -variable nonlinear root-finding sub-problem (56) to obtain  $\boldsymbol{\mu}^{(s)}$  with current  $\lambda_0$  and  $\lambda_1$  (see Remark 8)
5:   // Gradient Descent to Update  $\lambda_0$  and  $\lambda_1$ 
6:    $\lambda_0 \leftarrow \lambda_0 - \alpha \cdot \frac{\partial f(\boldsymbol{\lambda})}{\partial \lambda_0}$  // See (58)
7:    $\lambda_1 \leftarrow \lambda_1 - \alpha \cdot \frac{\partial f(\boldsymbol{\lambda})}{\partial \lambda_1}$  // See (59)
8:   // Projection
9:   if  $\lambda_0 < 0$  then  $\lambda_0 \leftarrow 0$ 
10:  end if
11:  // Next Iteration
12:   $s \leftarrow s + 1$ 
13:  // Stopping Rule
14:  if  $s > S$  or  $\text{abs}(\frac{\partial f(\boldsymbol{\lambda})}{\partial \lambda_1}) < \epsilon$  then
15:    if  $1 \neq \sum_i \mu_i^{(s)}$  then // Early Stopping Applied
16:       $\mu_i^{(s)} \leftarrow \mu_i^{(s)} / \sum_j \mu_j^{(s)}, \forall i \in [N]$ ,
17:    end if
18:    break while
19:  end if
20: end while
Output:  $\boldsymbol{\mu}^{(s)}$ 

```

---

**Remark 8:** We discuss the  $N$ -variate root-finding problem  $-2\mathbf{A}\boldsymbol{\mu} + \mathbf{b} - \lambda_1 \cdot \mathbf{1} + \lambda_0 \cdot (-\ln \boldsymbol{\mu} - \mathbf{1} + \ln \bar{\boldsymbol{\mu}}) = \mathbf{0}$  on  $\boldsymbol{\mu} \geq \mathbf{0}$ . Let  $\mathbf{g}(\boldsymbol{\mu}) := -2\mathbf{A}\boldsymbol{\mu} + \mathbf{b} - \lambda_1 \cdot \mathbf{1} + \lambda_0 \cdot (-\ln \boldsymbol{\mu} - \mathbf{1} + \ln \bar{\boldsymbol{\mu}})$ . One may verify that  $d\mathbf{g}(\boldsymbol{\mu})/d\boldsymbol{\mu} \prec \mathbf{0}$  (i.e.,  $\mathbf{g}$  is a monotonically decreasing function in  $\boldsymbol{\mu}$ ),  $\mathbf{g}(\mathbf{0}) \rightarrow \infty$ , and  $\mathbf{g}(\infty) \rightarrow -\infty$ . Therefore, at least one root of  $\mathbf{g}(\boldsymbol{\mu}) = \mathbf{0}$  exists and Newton's method can be used to find it.  $\square$

**Remark 9:** If the 2-norm constraint  $\|\boldsymbol{\mu} - \bar{\boldsymbol{\mu}}\|_2 \leq \theta_0$  is used to replace the KL divergence constraint, then the root-finding procedure would be significantly simplified. Therefore, in practice, to save computational time, one may choose the 2-norm constraint  $(\boldsymbol{\mu} - \bar{\boldsymbol{\mu}})^\top (\boldsymbol{\mu} - \bar{\boldsymbol{\mu}}) \leq \theta_0^2$ . Another choice to reduce the computational complexity is to use the Frank-Wolfe method (i.e., linearization of the objective function) as in Proposition 5.  $\square$

APPENDIX J  
THE STANDARD IMM FILTER

The implementation details of the interactive multiple model (IMM) method is given in Algorithm 3. The results in Step 2 (see Line 25) are due to (2) and (4) where  $\mu_{j,k|k-1}$  and  $\mu_{j,k|k}$  are prior and posterior model probabilities of the  $j^{\text{th}}$  model, respectively. The prior model probability, model likelihood, and posterior model probability of the  $j^{\text{th}}$  model are calculated in Step 1.5 (see Line 18), Step 1.6 (see Line 20), and Step 1.7 (see Line 22), respectively. See [3], [5] (in the reference list of the main body of the paper) for more information.

---

**Algorithm 3** Interactive Multiple Model Algorithm [3], [5]

---

**Definition:** Let  $\hat{\mathbf{x}}_{j,k|k-1}$  denote the prior state estimate provided by the  $j^{\text{th}}$  model and  $\mathbf{P}_{j,k|k-1}$  the corresponding state estimation error covariance. Let  $\hat{\mathbf{x}}_{j,k|k}$  denote the posterior state estimate provided by the  $j^{\text{th}}$  model and  $\mathbf{P}_{j,k|k}$  the corresponding state estimation error covariance; Let  $\hat{\mathbf{x}}_{k|k}$  denote the combined posterior state estimate of the  $N$  models and  $\mathbf{P}_{k|k}$  the corresponding state estimation error covariance; Let  $\mu_{j,k|k-1}$  and  $\mu_{j,k|k}$  be the prior and posterior model probability of the  $j^{\text{th}}$  model at the time  $k$ , respectively; Let  $\{\pi_{ij}\}_{i,j=1,2,\dots,N}$  be the model transition probability matrix.

**Initialization:**  $\forall j \in [N]$ , initialize  $\mu_{j,0|0}$ ,  $\hat{\mathbf{x}}_{j,0|0}$ , and  $\mathbf{P}_{j,0|0}$ .

**Remark:** In literature, prior and posterior state estimate are also known as predicted and updated state estimate, respectively.

**Input:**  $\mathbf{y}_k$ ,  $k = 1, 2, 3, \dots$

```

1: while true do
2:   // (Step 1) At Time  $k$ 
3:   for  $j = 1 : N$  do
4:     // (Step 1.1) Transition Probability From  $i^{\text{th}}$  Model at Time  $k-1$  To  $j^{\text{th}}$  Model at Time  $k$ 
5:      $\mu_{ij,k|k-1} = \frac{\pi_{ij} \cdot \mu_{i,k-1|k-1}}{\sum_{i=1}^N \pi_{ij} \cdot \mu_{i,k-1|k-1}}$ 
6:     // (Step 1.2) Initialize the  $j^{\text{th}}$  Filter
7:      $\hat{\mathbf{x}}_{j,k-1|k-1}^0 = \sum_{i=1}^N \mu_{ij,k|k-1} \cdot \hat{\mathbf{x}}_{i,k-1|k-1}$ 
8:      $\mathbf{P}_{j,k-1|k-1}^0 = \sum_{i=1}^N \mu_{ij,k|k-1} \cdot \left\{ \mathbf{P}_{i,k-1|k-1} + (\hat{\mathbf{x}}_{i,k-1|k-1} - \hat{\mathbf{x}}_{j,k-1|k-1}^0)(\hat{\mathbf{x}}_{i,k-1|k-1} - \hat{\mathbf{x}}_{j,k-1|k-1}^0)^{\top} \right\}$ 
9:     // (Step 1.3) Prior Estimation of the  $j^{\text{th}}$  Filter (i.e., Time Update)
10:     $\hat{\mathbf{x}}_{j,k|k-1} = \mathbf{F}_{j,k-1} \hat{\mathbf{x}}_{j,k-1|k-1}^0$ 
11:     $\mathbf{P}_{j,k|k-1} = \mathbf{F}_{j,k-1} \mathbf{P}_{j,k-1|k-1}^0 \mathbf{F}_{j,k-1}^{\top} + \mathbf{G}_{j,k-1} \mathbf{Q}_{j,k-1} \mathbf{G}_{j,k-1}^{\top}$ 
12:    // (Step 1.4) Posterior Estimation of the  $j^{\text{th}}$  Filter (i.e., Measurement Update)
13:     $\mathbf{r}_{j,k} = \mathbf{y}_k - \mathbf{H}_{j,k} \hat{\mathbf{x}}_{j,k|k-1}$  // Innovation
14:     $\mathbf{S}_{j,k} = \mathbf{H}_{j,k} \mathbf{P}_{j,k|k-1} \mathbf{H}_{j,k}^{\top} + \mathbf{R}_{j,k}$  // Innovation Covariance
15:     $\mathbf{K}_{j,k} = \mathbf{P}_{j,k|k-1} \mathbf{H}_{j,k}^{\top} \mathbf{S}_{j,k}^{-1}$  // Filter Gain
16:     $\hat{\mathbf{x}}_{j,k|k} = \hat{\mathbf{x}}_{j,k|k-1} + \mathbf{K}_{j,k} \cdot \mathbf{r}_{j,k} = \hat{\mathbf{x}}_{j,k|k-1} + \mathbf{P}_{j,k|k-1} \mathbf{H}_{j,k}^{\top} \mathbf{S}_{j,k}^{-1} \cdot [\mathbf{y}(k) - \mathbf{H}_{j,k} \hat{\mathbf{x}}_{j,k|k-1}]$ 
17:     $\mathbf{P}_{j,k|k} = \mathbf{P}_{j,k|k-1} - \mathbf{P}_{j,k|k-1} \mathbf{H}_{j,k}^{\top} \mathbf{S}_{j,k}^{-1} \mathbf{H}_{j,k} \mathbf{P}_{j,k|k-1}$ 
18:    // (Step 1.5) Prior Probability of the  $j^{\text{th}}$  Model
19:     $\mu_{j,k|k-1} = \sum_{i=1}^N \pi_{ij} \cdot \mu_{i,k-1|k-1}$ 
20:    // (Step 1.6) Likelihood of the  $j^{\text{th}}$  Model
21:     $\lambda_{j,k} = \mathcal{N}(\mathbf{r}_{j,k}; \mathbf{0}, \mathbf{S}_{j,k})$ 
22:    // (Step 1.7) Posterior Probability of the  $j^{\text{th}}$  Model
23:     $\mu_{j,k|k} = \frac{\mu_{j,k|k-1} \cdot \lambda_{j,k}}{\sum_{i=1}^N \mu_{i,k|k-1} \cdot \lambda_{i,k}}$ 
24:   end for
25:   // (Step 2) Combined Posterior State Estimate
26:    $\hat{\mathbf{x}}_{k|k} = \sum_{j=1}^N \mu_{j,k|k} \cdot \hat{\mathbf{x}}_{j,k|k}$ 
27:    $\mathbf{P}_{k|k} = \sum_{j=1}^N \mu_{j,k|k} \cdot \left\{ \mathbf{P}_{j,k|k} + (\hat{\mathbf{x}}_{j,k|k} - \hat{\mathbf{x}}_{k|k})(\hat{\mathbf{x}}_{j,k|k} - \hat{\mathbf{x}}_{k|k})^{\top} \right\}$ 
28:   // (Step 3) Next Time Step
29:    $k \leftarrow k + 1$ 
30: end while
Output:  $\hat{\mathbf{x}}_{k|k}, \mathbf{P}_{k|k}, \mu_{j,k|k}$ 

```

---