

Stochastic Observability and Filter Stability under Several Criteria

Curtis McDonald and Serdar Yüksel

Abstract—Despite being a foundational concept of modern systems theory, there have been few studies on observability of non-linear stochastic systems under partial observations. In this paper, we introduce a definition of observability for stochastic non-linear dynamical systems which involves an explicit functional characterization. To justify its operational use, we establish that this definition implies filter stability under mild continuity conditions: an incorrectly initialized non-linear filter is said to be stable if the filter eventually corrects itself with the arrival of new measurement information. Numerous examples are presented and a detailed comparison with the literature is reported. We also establish implications for various criteria for filter stability under several notions of convergence such as weak convergence, total variation, and relative entropy. These findings are connected to robustness and approximations in partially observed stochastic control.

Keywords: Observability, Non-Linear Filtering, Filter Stability, Merging

I. INTRODUCTION

Observability is one of the most important and foundational concepts of modern systems and control theory with implications at the heart of its theory and applications [10], [30], [38], [39]. For deterministic linear systems, observability is defined as the exact recovery of any initial condition with measurements available until some finite time, and is characterized by an observability rank condition in both continuous and discrete-time [11]. For linear systems, such an observability definition is global (as it applies for all initial states) and is also directly applicable to stochastic counterparts of deterministic linear systems. For non-linear systems, however, due to the challenges in the analysis which prevent globality, the analysis is significantly more nuanced both for deterministic and stochastic setups. See Section II-D for a detailed discussion.

We study the stochastic setup in this paper. Let us now introduce the probabilistic setup for a Hidden Markov Model (HMM) or Partially Observed Markov Process (POMP). Let $(\mathcal{X}, \mathcal{Y})$ be complete, separable and metric (Polish) spaces equipped with their Borel sigma fields $\mathcal{B}(\mathcal{X})$ and $\mathcal{B}(\mathcal{Y})$. \mathcal{X} will be called the state space, and \mathcal{Y} the measurement space. Let $\mathcal{P}(\mathcal{X})$ and $\mathcal{P}(\mathcal{Y})$ be the set of probability measures on these spaces. Define the transition kernel T and measurement channel G as the mappings

$$T : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{X}) \quad G : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{Y})$$

C. McDonald is with the Department of Statistics and Data Science at Yale University, United States of America (e-mail: curtis.mcdonald@yale.edu). S. Yüksel is with the Department of Mathematics and Statistics, Queen's University, Kingston, Ontario, Canada, K7L 3N6 (e-mail: yuksel@mast.queensu.ca). Research is supported by the Natural Sciences and Engineering Research Council of Canada. Some preliminary results of this submission were presented at the 2018 Annual Allerton Conference.

$$x \mapsto T(dx'|x) \quad x \mapsto G(dy|x)$$

The system is initialized with a state $X_0 \in \mathcal{X}$ drawn from a prior measure μ on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$. The state is then randomly updated via the transition kernel T which makes the state process $\{X_n\}_{n=0}^{\infty}$ a Markov Chain with initial measure μ and transition kernel T .

However, the state is not available at the observer, instead at time n the observer sees the observation Y_n where the conditional distribution of $Y_n|X_n$ is determined by the measurement channel G .

By stochastic realization arguments [23, Lemma 1.2], [9, Lemma 3.1], we can also view an equivalent construction of the system dynamics. Let $\{Z_n\}_{n=0}^{\infty}$ and $\{W_n\}_{n=0}^{\infty}$ be independent identically distributed (i.i.d) \mathcal{Z} -valued noise processes, where \mathcal{Z} can be taken to be $[0, 1]$ or \mathbb{R} (or any other Polish space), without any loss of generality. Consider a partially observed dynamical system with the following model.

$$\begin{aligned} X_{n+1} &= b(X_n, W_n) \\ Y_n &= h(X_n, Z_n), \end{aligned} \quad (1) \quad (2)$$

where W_n and Z_n can be assumed to take values from $[0, 1]$ or \mathbb{R} . Here b defines the system dynamics and defines a transition kernel T for the Markov chain X_n . Assuming Z_n has measure Q in \mathcal{Z} , the measurement function h defines the measurement channel G which is the pushforward measure of Q under $h(x, \cdot)$. Throughout the paper we will work with either the general kernel and measurement channel notation T, G or with the specific functional form using b, h when convenient.

Thus, the observer needs to compute the conditional probability on the hidden variable X_n using the information available up to time $n \in \mathbb{Z}_+$. We have that $\{X_n, Y_n\}_{n=0}^{\infty}$ is a Markov chain, and we will denote P^μ as the probability measure on $\Omega = \mathcal{X}^{\mathbb{Z}_+} \times \mathcal{Y}^{\mathbb{Z}_+}$, endowed with the product topology, (and thus $\omega \in \Omega$ is a sequence of states and measurements $\omega = \{(x_i, y_i)\}_{i=0}^{\infty}$) where $X_0 \sim \mu$. Such a stochastic system is referred to as a Partially Observed Markov Process (POMP) (also called Hidden Markov Model) throughout the paper.

Definition I.1. We define the one step predictor as the sequence

$$\pi_{n-}^\mu(\cdot) = P^\mu(X_n \in \cdot | Y_0, \dots, Y_{n-1}), \quad n \in \mathbb{Z}_+$$

and we define the non-linear filter as the sequence

$$\pi_n^\mu(\cdot) = P^\mu(X_n \in \cdot | Y_0, \dots, Y_n), \quad n \in \mathbb{Z}_+.$$

Both of the above are regular conditional probability sequences defined on \mathcal{X} . We will use the notation $Y_{[0,n]} =$

Y_0, \dots, Y_n to represent finite sets of random variables, and $Y_{[0,\infty)} = Y_0, Y_1, \dots$ to represent infinite sequences. The recursive update equations for the filter or the predictor are known as the non-linear filtering equations. Let us, for the time being, assume the existence of a likelihood function $g(x, y)$ for the measurement channel defined as follows. The measurement channel G is called dominated if there exists a reference measure λ such that $\forall x \in \mathcal{X}, G(Y_n \in \cdot | X_n = x) \ll \lambda$ where the notation " \ll " denotes absolute continuity. We can then utilize a likelihood function $g(x, y) = \frac{dG(Y_n \in \cdot | X_n = x)}{d\lambda}(y)$ and write the filter π_{n+1}^μ recursively in terms of π_n^μ and $Y_{n+1} = y_{n+1}$ explicitly as a Bayesian update:

$$\begin{aligned} \pi_{n+1}^\mu(dx_{n+1}) &= F(\pi_n^\mu, y_{n+1})(dx_{n+1}) \\ &:= \frac{g(x_{n+1}, y_{n+1}) \int_{\mathcal{X}} T(dx_{n+1} | X_n = x) \pi_n^\mu(dx)}{\int_{\mathcal{X}} g(x_{n+1}, y_{n+1}) \int_{\mathcal{X}} T(dx_{n+1} | X_n = x) \pi_n^\mu(dx)} \end{aligned} \quad (3)$$

Suppose that an observer runs a non-linear filter assuming that the initial prior is ν , when in reality the prior distribution is μ . The observer receives the measurements and computes the filter π_n^ν for each n , but the measurement process is generated according to the true measure μ .

The operational question for observability is that of *filter stability*, namely, if we have two different initial probability measures μ and ν , when do we have that the filter processes π_n^μ and π_n^ν merge in some appropriate sense as $n \rightarrow \infty$. In essence, when will our observations Y_n be informative enough to correct our incorrect prior ν and result in an accurate conditional measure for the hidden state.

In Section I-A below, notations and definitions are presented. In Section II we present our main results. We present a detailed literature review after the statement of our main results in Section II-D. Examples of observable systems are given in Section III. Proofs are provided in Section IV.

A. Notation and preliminaries

Let $C_b(\mathcal{X})$ represent the set of continuous and bounded functions from $\mathcal{X} \rightarrow \mathbb{R}$.

Definition I.2. *Two sequences of probability measures P_n, Q_n merge weakly if $\forall f \in C_b(\mathcal{X})$ we have $\lim_{n \rightarrow \infty} |\int f dP_n - \int f dQ_n| = 0$.*

Definition I.3. *For two probability measures P and Q the total variation norm is defined as $\|P - Q\|_{TV} = \sup_{\|f\|_\infty \leq 1} |\int f dP - \int f dQ|$ where f is assumed measurable.*

Note that merging in total variation implies weak merging since $C_b(\mathcal{X})$ is a subset of the set of measurable and bounded functions. We also utilize the relative entropy (Kullback-Leibler divergence) between two probability measures, although it is not a metric (since it is not symmetric).

Definition I.4.

(i) For two probability measures P and Q we define the relative entropy as $D(P||Q) = \int \log \frac{dP}{dQ} dP = \int \frac{dP}{dQ} \log \frac{dP}{dQ} dQ$ where $P \ll Q$ and $\frac{dP}{dQ}$ denotes the Radon-Nikodym derivative of P with respect to Q .

(ii) Let X and Y be two random variables, let P and Q be two different joint measures for (X, Y) with $P \ll Q$. We define the (conditional) relative entropy between $P(X|Y)$ and $Q(X|Y)$ as

$$\begin{aligned} D(P(X|Y)||Q(X|Y)) &= \int \log \left(\frac{dP_{X|Y}}{dQ_{X|Y}}(x, y) \right) P(dx, y) \\ &= \int \left(\int \log \left(\frac{dP_{X|Y}}{dQ_{X|Y}}(x, y) \right) P(dx|Y = y) \right) P(dy) \end{aligned} \quad (4)$$

Some notational discussion is in order. For some probability measures such as $P^\mu(Y_{[0,n]} \in \cdot)$ or $P^\mu(X_n \in \cdot)$, it will be convenient to denote the random variable inside the measure and take out the set argument. When we take the relative entropy of such measures, to make the notation shorter, we will drop the " $\in \cdot$ " argument and write $D(P^\mu(Y_{[0,n]})||P^\nu(Y_{[0,n]}))$.

Note that in a conditional relative entropy, we are integrating the logarithm of the Radon-Nikodym derivative of the conditional measures $P(X|Y)$ and $Q(X|Y)$ over the joint measure of P on (X, Y) . The second equality (4) shows that this can be thought of as the expectation of the relative entropy $D(P(X|Y = y)||Q(X|Y = y))$ at specific realizations of $Y = y$, where the expectation is over the marginal measure of P on Y . When we apply this to the filter, π_n^μ and π_n^ν are realizations of the filter for specific measurements, therefore when we discuss their relative entropy, we take the expectation over the marginal of P^μ on $Y_{[0,n]}$. We write this as $E^\mu[D(\pi_n^\mu||\pi_n^\nu)]$ where $D(\pi_n^\mu||\pi_n^\nu)$ plays the role of the inner integral in (4).

We now introduce some additional notation that will be useful when dealing with sigma fields rather than random variables directly. Strictly speaking, we have two probability measures P^μ and P^ν on $(\mathcal{X}^{\mathbb{Z}^+} \times \mathcal{Y}^{\mathbb{Z}^+}, \mathcal{B}(\mathcal{X}^{\mathbb{Z}^+} \times \mathcal{Y}^{\mathbb{Z}^+}))$. We denote by $\mathcal{F}_{a,b}^{\mathcal{X}}$ the sigma field generated by (X_a, \dots, X_b) and similarly for \mathcal{Y} . We also write $\mathcal{F}_n^{\mathcal{X}}$ for the sigma field generated by X_n . We then have $\mathcal{F}_{0,\infty}^{\mathcal{X}} \vee \mathcal{F}_{0,\infty}^{\mathcal{Y}}$ as the sigma field generated by all state and measurement sequences. When we write $P^\mu(X_{[0,n]})$ we are discussing the measure P^μ restricted to the sigma field $\mathcal{F}_{0,n}^{\mathcal{X}}$ which we will denote $P^\mu|_{\mathcal{F}_{0,n}^{\mathcal{X}}}$. Similarly for some set $A \in \mathcal{F}_{0,\infty}^{\mathcal{X}} \vee \mathcal{F}_{0,\infty}^{\mathcal{Y}}$ we write $P^\mu((X_{[0,\infty)}, Y_{[0,\infty)}) \in A|Y_{[0,n]})$ as the conditional measure of P^μ with respect to the sigma field $\mathcal{F}_{0,n}^{\mathcal{Y}}$, which we denote $P^\mu|_{\mathcal{F}_{0,n}^{\mathcal{Y}}}$. We can also consider restricting and conditioning simultaneously, this for example is the case with the non-linear filter: $\pi_n^\mu(\cdot) = P^\mu(X_n \in \cdot | Y_{[0,n]}) = P^\mu|_{\mathcal{F}_n^{\mathcal{X}}}|_{\mathcal{F}_{0,n}^{\mathcal{Y}}}$. The key relationship between relative entropy and total variation is Pinsker's inequality (see e.g., [16]) which states that for two probability measures P and Q we have that $\|P - Q\|_{TV} \leq \sqrt{\frac{2}{\log(e)} D(P||Q)}$.

Criteria for stability. We note the following definitions for filters, but they can also be defined for predictors.

Definition I.5. (i) A filter process is stable in the sense of weak merging in expectation if for any $f \in C_b(\mathcal{X})$ and any prior ν with $\mu \ll \nu$ we have

$$\lim_{n \rightarrow \infty} E^\mu \left[\left| \int f d\pi_n^\mu - \int f d\pi_n^\nu \right| \right] = 0.$$

(ii) A filter process is stable in the sense of weak merging P^μ almost surely (a.s.) if there exists a set of measurement sequences $A \subset \mathcal{Y}^{\mathbb{Z}^+}$ with P^μ probability 1 such that for any sequence in A , for any $f \in C_b(\mathcal{X})$ and any prior ν with $\mu \ll \nu$ we have

$$\lim_{n \rightarrow \infty} \left| \int f d\pi_n^\mu - \int f d\pi_n^\nu \right| = 0.$$

(iii) A filter process is stable in the sense of total variation in expectation if for any measure ν with $\mu \ll \nu$ we have

$$\lim_{n \rightarrow \infty} E^\mu [\|\pi_n^\mu - \pi_n^\nu\|_{TV}] = 0.$$

(iv) A filter process is stable in the sense of total variation P^μ a.s. if for any measure ν with $\mu \ll \nu$ we have

$$\lim_{n \rightarrow \infty} \|\pi_n^\mu - \pi_n^\nu\|_{TV} = 0 \quad P^\mu \text{ a.s.}$$

(v) A filter process is stable in relative entropy if for any measure ν with $\mu \ll \nu$:

$$\lim_{n \rightarrow \infty} E^\mu [D(\pi_n^\mu \|\pi_n^\nu)] = 0.$$

(vi) For $f : \mathcal{X} \rightarrow \mathbb{R}$ define the Lipschitz norm

$$\|f\|_L = \sup \left\{ \frac{|f(x) - f(y)|}{d(x, y)} \mid d(x, y) \neq 0 \right\}$$

With $BLip := \{f : \|f\|_L \leq 1, \|f\|_\infty \leq 1\} \subset C_b(\mathcal{X})$ we define the bounded Lipschitz (BL) metric as

$$\|P - Q\|_{BL} = \sup_{f \in BLip} \left| \int f dP - \int f dQ \right|.$$

A system is then stable in the sense of BL-merging P^μ a.s. if we have $\|\pi_n^\mu - \pi_n^\nu\|_{BL} \rightarrow 0 \quad P^\mu \text{ a.s.}$

We note that *merging* of probability measures is different from the *convergence* of a sequence of probability measures to a limit measure. In convergence, we have some sequence P_n and a static limit measure P ; in merging we have two sequences P_n and Q_n which may not individually have limits, but come closer together for large n in one of the merging notions defined previously [17]

II. STATEMENT OF THE MAIN RESULTS AND LITERATURE REVIEW

A. Stochastic non-linear observability

We first introduce our notion of an observable system.

Definition II.1. (i) [One Step Observability] A POMP is said to be one step observable if for every $f \in C_b(\mathcal{X})$, $\epsilon > 0$, \exists a measurable and bounded function $g : \mathcal{Y} \rightarrow \mathbb{R}$ such that

$$\left\| f(\cdot) - \int_{\mathcal{Y}} g(y) G(dy|\cdot) \right\|_\infty < \epsilon$$

(ii) [N Step Observability] A POMP is said to be N step observable if for every $f \in C_b(\mathcal{X})$, $\epsilon > 0$, \exists a measurable and bounded function $g : \mathcal{Y}^N \rightarrow \mathbb{R}$ such that

$$\left\| f(\cdot) - \int_{\mathcal{Y}^N} g(y_{[1, N]}) P(dy_{[1, N]} | X_1 = \cdot) \right\|_\infty < \epsilon, \quad (5)$$

where we note that the conditional probability $P(dy_{[1, N]} | X_1 = x_1)$ is independent of the prior measure.

(iii) [Observability] A POMP is observable if for every $f \in C_b(\mathcal{X})$ and every $\epsilon > 0$ there exist $N \in \mathbb{N}$ and a measurable and bounded function g (both possibly dependent on f and ϵ) such that (5) applies. Note that if a POMP is N step observable for some finite $N \in \mathbb{N}$ then it is observable, but the reverse implication is not necessarily the case.

A number of remarks are in order.

Remark: II.2. In the definition above, we can instead of $C_b(\mathcal{X})$ consider any dense subset in $C_b(\mathcal{X})$. For example, if \mathcal{X} is a compact subset of \mathbb{R}^k , we can consider polynomials as these form a dense subset, or we can consider smooth functions defined on \mathcal{X} , or functions which are expressed as linear combinations of harmonics, Haar wavelets etc. An example is provided in Section III-B.

Remark: II.3. [One-step observability and universality in the controlled setup] The definition of one step observability is a specific case of N step observability, however the distinction is useful for at least two reasons: (i) One-step observability is often easier to check since one doesn't need to consider the effect of the state transition kernel, as this definition is only concerned with the measurement channel itself. On the other hand, there exist many setups where a system is observable, but not one step observable; see e.g. Section III-A. (ii) Even though in this paper we consider a control-free setup, in a controlled context studied in a companion paper [43], it follows that one step observability would be independent of any control policy (that is, observability would be universal over all policies and associated filter stability results apply under any control policy), but $N > 1$ step observability would be handled much more cautiously as this condition would be dependent on the control policy adopted. Recently, filter stability results have been shown to be consequential in showing near-optimality of finite memory control policies and associated learning theoretic results for Partially Observed Markov Decision Processes [33], [34]. Accordingly, one-step observability results are particularly applicable for such scenarios.

Remark: II.4. If the measurement kernel satisfies an absolute continuity condition so that $G(dy|x) = h(x, y)\lambda(dy)$ and if there exists a finite measure K such that $G(dy|x) \leq K(dy)$ (so that the family of kernels $\{G(dy|x), x \in \mathcal{X}\}$ is majorized by K leading to a uniformly countable additive family of measures), then by Lusin's theorem [18, Theorem 7.5.2] and the extension theorem of Tietze [19, Theorem 4.1], we can replace g in the above with a continuous function g_c . The relaxation to such continuous g_c is useful when one would like to approximate the channels with those that are quantized. This then leads to an easier way to test observability via a rank condition, e.g. when \mathcal{X} is finite; see Section III-A.

Remark: II.5. One should note that the definition is not one of invertibility; it only requires that there exists some g and

N such that the error between the conditional expectation of $g(y_{[1,N]})$ given $X_1 = x$, and $f(x)$ is small. In particular, X_1 is not necessarily, even approximately, recoverable given the measurements. Invertibility, however, would be a special case being a sufficient condition, as we will see in the examples.

Remark: II.6. [Recovery of Initial Probability Measure] By our definition of observability, for every $f \in C_b(\mathbb{X})$, the value is $\langle \mu, f \rangle := E_\mu[f(X)]$ determinable with arbitrary precision by the measurements (since f is recovered, uniformly over a given compact set, with arbitrary precision). Since a countable collection of continuous and bounded functions uniquely distinguish probability measures [20, Theorem 3.4.5] (that is, such continuous and bounded functions form a separating class, see also [6, p. 13]), this amounts to the recovery of the initial probability measure as more and more measurements are collected. This then leads to the conclusion that our definition implies Van Handel's definition given in (7), noted further below (note that this also applies for non-compact setups under Definition II.14 as every individual probability measure is tight).

B. Filter stability under the observability definition

The presented observability definition leads to predictor stability in the following sense.

Theorem II.7. *Let*

$$P^\mu|_{\mathcal{F}_{0,\infty}^y} \ll P^\nu|_{\mathcal{F}_{0,\infty}^y}. \quad (6)$$

If the POMP is observable, then π_{n-}^μ and π_{n-}^ν merge weakly as $n \rightarrow \infty$, P^μ a.s.

A sufficient condition for (6) is that the priors satisfy $\mu \ll \nu$. The following assumption will allow us to use the recent results in [31] (see also [21]) and conclude the weak merging of the filter in expectation from the almost sure convergence of the predictor.

Assumption II.8. *The measurement channel is continuous in total variation. That is for any sequence a_n with $\lim_{n \rightarrow \infty} a_n = a \in \mathcal{X}$ we have $\|G(\cdot|X_0 = a_n) - G(\cdot|X_0 = a)\|_{TV} \rightarrow 0$.*

An example of a channel which is continuous in total variation is as follows [31, Section 2.2]: $Y_n = F(X_n) + W_n$ where F is continuous and W_n admits a continuous density function (such as a Gaussian), where an analysis based on convolution and Scheffé's lemma leads to the conclusion.

Theorem II.9. *Let Assumption II.8 hold, if the predictor merges weakly P^μ a.s., then the filter merges weakly in expectation.*

1) *Localized Observability for Non-Compact Signal Spaces:* While the definition of observability that we introduced is valid for both compact and non-compact state spaces, it may be difficult to satisfy the definition in a non-compact state space with a uniform bound on the approximation error. This will be relaxed in the following, where we assume \mathcal{X} to be Euclidean.

Definition II.10. *Given a compact set K , a POMP is called K locally predictable if there exists a sequence of $\mathcal{F}_{0,n-1}^y$ (with $n \in \mathbb{N}$) measurable mappings (random variables) $a_n : \mathcal{Y}^n \rightarrow \mathcal{X}$ such that*

$$\pi_{n-}^\nu(K + a_n) = 1 \quad P^\mu \text{ a.s.}$$

for every $\mu \ll \nu$.

This definition can be interpreted as follows. Regardless of the prior ν , upon seeing observations $Y_{[0,n-1]}$ we can be sure X_n lives in a compact set $K_n = K + a_n$. We can think of a_n as a ‘‘centring’’ value based on the observations $Y_{[0,n-1]}$ and K as the compact set around this centring value in which X_n must live conditioned on observations $Y_{[0,n-1]}$. This is then paired with a definition of local observability.

Definition II.11. *Given a compact set K , a POMP is called K locally observable if for every continuous and bounded function f , every sequence of numbers a_n , and every $\epsilon > 0$, there exists a sequence of uniformly bounded measurable functions g_n such that*

$$\sup_{x \in K + a_n} \left| f(x) - \int_{\mathcal{Y}} g_n(y) G(dy|x) \right| \leq \epsilon$$

for every $n \in \mathbb{N}$.

Theorem II.12. *Assume $\mu \ll \nu$ and that there exists a compact set K such that the POMP is K locally predictable and K locally observable. Then the predictor merges weakly P^μ a.s..*

The result above is intuitive as we can specify the compact set K , and the shifted sets $K_n = K + a_n$, over which we must approximate the function. However, the definitions can also be constructed taking a more relaxed approach and appealing to tightness rather than a probability one statement, but in this case it is more difficult to satisfy the local definition of observability.

Definition II.13. *A POMP is called locally predictable if there exists a sequence of $\mathcal{F}_{0,n-1}^y$ (with $n \in \mathbb{N}$) measurable mappings $a_n : \mathcal{Y}^n \rightarrow \mathcal{X}$ such that the family of measures*

$$\tilde{\pi}_{n-}^\nu(\cdot) := \pi_{n-}^\nu(\cdot + a_n)$$

for every $\mu \ll \nu$, is a uniformly tight family of measures.

Definition II.14. *A POMP is called locally observable if for every continuous and bounded function f , every compact set K , every sequence of numbers a_n , and every $\epsilon > 0$, there exists a sequence of uniformly bounded measurable functions g_n such that*

$$\sup_{x \in K + a_n} \left| f(x) - \int_{\mathcal{Y}} g_n(y) G(dy|x) \right| \leq \epsilon,$$

$$\sup_{x \notin K + a_n} \left| \int_{\mathcal{Y}} g_n(y) G(dy|x) \right| \leq 2\|f\|_\infty$$

for every $n \in \mathbb{N}$.

Theorem II.15. *Assume $\mu \ll \nu$ and that the POMP is locally predictable and locally observable. Then the predictor merges weakly P^μ a.s..*

An example is in Section III-D.

C. Relations between various criteria for filter stability

It follows from Pinsker's inequality that relative entropy merging implies total variation merging, which in turn implies weak merging (by Definitions I.2 and I.3). In this section we are interested conditions for when the converse direction holds, i.e. weak merging implies total variation or relative entropy merging. Recall the definition that for the measurement channel $G(Y_n \in \cdot | X_n = x)$ to be *dominated* in the sense that there exists a reference measure λ such that $\forall x \in \mathcal{X}, G(dy|x_n = x) \ll \lambda$. Then, we define the Radon-Nikodym derivative

$$g(x, y) := \frac{dG(y_n \in \cdot | x_n = x)}{d\lambda(\cdot)}(y)$$

which serves as a likelihood function (a conditional probability density function).

- Assumption II.16.** (i) $T(\cdot|x)$ is absolutely continuous with respect to a dominating measure ϕ for every $x \in \mathcal{X}$, so that $t(x_1, x) = \frac{dT(\cdot|x)}{d\phi}(x_1)$ where t is continuous in x for every $x_1 \in \mathcal{X}$.
(ii) $g(x, y)$ is bounded and continuous in x for every fixed y . Furthermore, $g(x, y) > 0$ for all $x \in \mathcal{X}, y \in \mathcal{Y}$.

Assumption II.17. $T(\cdot|x)$ is absolutely continuous with respect to a dominating measure ϕ for every $x \in \mathcal{X}$, so that $t(x_1, x) = \frac{dT(\cdot|x)}{d\phi}(x_1)$. The family of (conditional densities) $\{t(\cdot, x)\}_{x \in \mathcal{X}}$ is uniformly bounded and equicontinuous.

Theorem II.18. Let $\mu \ll \nu$. Let either one of the following hold:

- i) Assumption II.16, or
- ii) Assumption II.17.

Then, if the predictor is stable in the weak sense P^μ a.s. then it is also stable in total variation P^μ a.s..

Since the total variation of any two probability measures is uniformly bounded, stability in the almost sure sense implies that in expectation (with the same also holding for predictors). Thus, Theorem II.18 also presents condition for predictor merging in total variation in expectation.

Theorem II.19. The filter merges in total variation in expectation if and only if the predictor merges in total variation in expectation.

Recently, the filter stability results under total variation in expectation have been shown to be consequential in showing the optimality of finite memory control policies in Partially Observed Markov Decision Processes (see [34, Section 4.3 and Theorem 9] and [33, Theorems 3.2, 3.3 and 4.1]).

Theorem II.20. Assume there exists some finite n such that $E^\mu[D(\pi_n^\mu || \pi_n^\nu)] < \infty$ and some m such that $D(P^\mu|_{\mathcal{F}_{0,m}^\nu} || P^\nu|_{\mathcal{F}_{0,m}^\nu}) < \infty$. Then the filter is stable in relative entropy if and only if it is stable in total variation in expectation.

We note that both of the conditions on the finiteness of relative entropies in Theorem II.20 are minor and hold for example

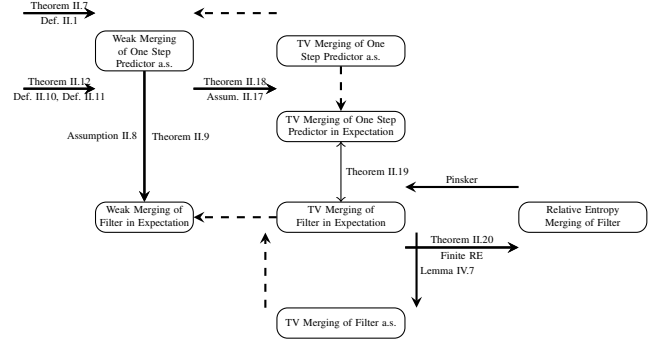


Fig. 1: Flow of Ideas and Conditions for Filter Stability

if $D(\mu||\nu) < \infty$. In the special setup where the measurement sigma field is the trivial one (with no information), or more generally Y_n is independent of X_n ; the above recovers the following result which generalizes Barron [3], [4] and Fritz [22], who had established the relative entropy convergence of sequences of probability measures for each time stage to the invariant measure for reversible Markov chains. This result also generalizes Theorem 5 of Harremoës and Holst [26] which considers countable state space chains with a uniform irreducibility assumption.

Theorem II.21. Let X_t be a Markov chain with π being its unique invariant probability measure. Let π_t denote the measure $P^{\pi_0}(X_t \in \cdot)$, where $X_0 \sim \pi_0$. Let $\pi_t \rightarrow \pi$ in total variation. If $D(\pi_{t_0}||\pi) < \infty$ for some $t_0 < \infty$, then

$$D(\pi_t||\pi) \downarrow 0.$$

In particular, for an aperiodic positive Harris recurrent Markov chain, if $D(\pi_{t_0}||\pi) < \infty$ for some $t_0 < \infty$, then $D(\pi_t||\pi) \downarrow 0$.

Proof. The proof follows directly from Theorem II.20. In the special case of positive Harris recurrence, the result follows since for aperiodic positive Harris recurrent Markov chains, $\pi_t \rightarrow \pi$ in total variation (see Theorem 13.0.1 in [44]). \square

D. Literature review and comparison of results

For deterministic linear systems, exact recovery of any initial condition with measurements available until some finite time is defined as observability and is characterized by an observability rank condition in both continuous and discrete-time [11]. For linear systems, such an observability definition is global (as it applies for all initial states) and is universal in the control policies applied, as the control policy does not affect the estimation errors (known as the *no-dual effect* [2] property). For non-linear systems, however, due to the challenges in the analysis which prevent globality as well as control-dependence, more modest and localized definitions are to be imposed: For deterministic continuous-time non-linear systems [28] and [49] present local indistinguishability conditions with subtle differences, and establish relations with Lie-theoretic characterizations which generalize observability rank conditions for non-linear systems defined locally. For discrete-time deterministic models, observability has also been

defined by invertibility or exact recovery of an initial state, locally, given measurements with finitely many observations. Nijmeier [45] developed discrete-time analogues of the observability notions presented in [28] (see also [49] for sampled continuous-time systems). Liu and Bitmead [41], [42] introduce a non-linear stochastic observability definition through entropy, where the conditional entropy of the hidden state given measurements not being the same as the unconditional entropy implies observability. Ugrinoovski [51] also presents an information theoretic formulation, and defines observability as an informativeness condition.

In the filtering literature for control systems, the classical setup involves the linear Gaussian system. The filter in this case is the celebrated Kalman filter, where the finite-dimensional Kalman filter is computed recursively using the Riccati equation. Under linear observability and controllability conditions, the Riccati equation admits a unique solution [10], [38], [39], which is the unique limit of the Riccati recursions regardless of the initialization. Thus, the Kalman Filter is stable with respect to incorrect, though still Gaussian, priors under the aforementioned conditions (we note that partial convergence and robustness results on the asymptotic equivalence of conditional expectations and linear estimates for non-Gaussian priors for linear systems are reported in [50]). The time-varying linear system setup has been studied in [1].

In a recent paper, we studied the implications on filter stability in robust control [43]. Implications on finite memory approximations in optimal stochastic control have been presented in [33], [34]. It is worth pointing out that there has been a recurrent theme on the duality between controllability and observability, for a recent work in this direction see [35], [36]. Filter stability for deterministic systems under noisy measurements has recently been studied in [47].

A strict version of our observability definition is captured in [13, Equation 1.7]. The idea there is to express, exactly, a continuous function $f(x)$ by integrating a measurable function $g(y)$ over the conditional distribution for Y given $X = x$. A fundamental result which pairs with observability is that of Blackwell and Dubins [7], an implication of which [13] independently arrived at. Blackwell and Dubins use martingale convergence theorem to show that if P and Q are two measures on a fully observed stochastic process $\{X_n\}_{n=0}^\infty$ with $P \ll Q$, then the conditional distributions on the future based on the past merge in total variation P a.s., that is P a.s.

$$\|P(X_{[n+1,\infty)} \in \cdot | X_{[0,n]}) - Q(X_{[n+1,\infty)} \in \cdot | X_{[0,n]})\|_{TV} \rightarrow 0.$$

[53] introduces a definition of observability for POMP. Namely, a system is observable if every prior results in a unique probability measure on the measurement sequences;

$$P^\mu|_{\mathcal{F}_{0,\infty}^y} = P^\nu|_{\mathcal{F}_{0,\infty}^y} \implies \mu = \nu. \quad (7)$$

[53] shows that the above leads to filter stability for continuous-time models with compact state space. [55] extends these results to non-compact state spaces, where *uniform observability* is introduced. The result of Blackwell and Dubins [7] is utilized to show that uniform observability would imply filter stability in bounded Lipschitz distance [52]. Nonetheless,

this condition is implicit; [52] only studies the measurement channel where $h(x, z) = f(x) + z$ where f^{-1} is uniformly continuous and Z must have an everywhere non-zero characteristic function (e.g. a Gaussian distribution). For a compact state space, [55] has established that uniform observability and observability are equivalent notions. We also note that for a finite state space with a non-degenerate measurement channel (i.e. likelihood function $g(x, y) > 0$), stability can be fully characterized via observability and a detectability condition [53], [56, Theorem V.2] or [12, Theorems 2.7 and 3.1].

As noted in Remark II.6, our definition implies (7). (7) is a statement of invertibility with no clear guidance on how to test this property; our definition is explicitly given in a test function formulation, making it more interpretable and easier to apply to various systems of interest. Additionally, in the work studied here, we consider discrete time processes and thus the predictor and the filter are distinct objects. Our definition of observability only implies the weak merging of the predictor almost surely, not the filter directly. Conditions are needed to relate the merging of the predictor to that of the filter. This is also addressed in our paper building also on recent results on the regularity properties of non-linear filters from [31] (see also [21]).

In early work by Kunita, [37], the stability of the filter process is studied in light of the limit sigma fields of the processes (e.g. $\mathcal{F}_{0,\infty}^y, \mathcal{F}_{0,\infty}^x$). Kunita's work unfortunately made a technical error on the exchange of orders in supremum and intersection operations on sigma fields: A concise derivation of the corrected result is presented in [14, Equation 1.10]: here, we are presented with a sufficient and necessary condition for the merging of the filter in total variation in expectation based on comparing the sigma fields $\mathcal{F}_{0,\infty}^y$ and $\bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^x \vee \mathcal{F}_{0,\infty}^y$. That is the filter merges in total variation in expectation if and only if:

$$E^\nu\left[\frac{d\mu}{d\nu}(X_0) | \mathcal{F}_{0,\infty}^y\right] = E^\nu\left[\frac{d\mu}{d\nu}(X_0) | \bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^x \vee \mathcal{F}_{0,\infty}^y\right] \quad P^\mu \text{ a.s.}$$

Relative entropy as a measure of discrepancy between the true filter and the incorrectly initialized filter is studied by Clark, Ocone, and Coumarbatch in [15]. Here, the authors considered the filtering problem in continuous time and with a dominated measurement channel. The authors established the relative entropy of the true filter and the incorrect filter as a supermartingale, and its convergence to a limit. However, the paper did not establish the convergence to zero. A notable setup where actual convergence (of the relative entropy) to zero is established is the (rather specific) Beneš filter studied in [46]. This problem also has relations to the relative entropy convergence of Markov chains to invariance: in the case where the measurements are trivial, the convergence problem reduces to what has been studied in [3], [4], [22], [26], [27] on relative entropy convergence of Markov chains to invariant measures.

Contributions and comparison with the literature. In view of the review above, our contributions are as follows:

- i) [Stochastic observability] In Section II-A, we present a definition of stochastic observability. This definition is functionally explicit and testable, and due to its functional

approximation characterization, it allows various analytical methods to be applicable for verification (see Remarks II.2 through Remark II.5).

Under this definition, we establish predictor stability (in the weak convergence/merging sense). We note that observability, for the discrete time case as studied here, only implies weak merging of the predictor almost surely, not the filter directly. This is addressed in our paper building also on recent results on the regularity properties of non-linear filters from [31] (see also [21]).

We also note that the Blackwell and Dubins theory of merging on which our approach builds (similar to [53], [55]), applies for infinite sequences of future events (i.e. $P^\mu(Y_{[n+1,\infty)}|Y_{[0,n]})$), this is utilized in our definition of N step observability leading to application examples of broad generality. Accordingly, our definition is not only a function of the measurement channel, but also of the system dynamics; unlike some related results in the literature.

Additionally, we provide several examples in Section III.

- ii) [On various convergence and merging criteria] We establish new results relating various criteria for filter stability (as depicted in Figure 1), independent of the mechanism used to arrive at filter stability: we study filter stability under weak merging and total variation merging in expectation and almost surely, as well as relative entropy. In Section II-C, (a) We place mild assumptions on the transition/measurement kernels to extend weak merging of the predictor to total variation merging. (b) We show that total variation merging of the predictor and filter are equivalent, and (c) under a mild finiteness condition on the relative entropy sequence, we also establish equivalence of relative entropy merging and total variation merging.

Using the chain rule for relative entropy, the relative entropy error was shown to be a non-increasing sequence by Clark et.al. [15], but *its convergence to zero was not established*, except for the specific case of the Beneš filter in [46]. Theorem II.20 establishes the equivalence between relative entropy merging and total variation, and thus convergence of the relative entropy error to zero is proven here (we note that this is a result which is hinted at in the literature, see [14, Remark 4.2] or [54, Remark 5.9] but not explicitly proven). This result applies to setups beyond filter stability: In the case where the measurements are trivial, Theorem II.21 generalizes [3], [4], [22], [26] on relative entropy convergence of Markov chains to invariant measures, where the first references due to Barron and Fritz had considered reversible Markov chains and the latter due to Harremoës and Holst focused on countable state Markov chains under a uniform irreducibility assumption. On Theorem II.19, we note first that much of the literature focuses on continuous time, where the predictor is not used in the analysis. In discrete time, [52, Lemma 4.2] proves that the merging of the predictor in total variation in expectation implies that of the filter. However this result relies on a domination assumption in the measurement channel and the specific structure of the filter recursion equation [14, Equation 1.4]. Theorem II.19 is, accordingly, a more general result.

- iii) [Implications to near optimality of finite window policies

in POMDPs] Our findings lead to practically relevant and mathematically consequential implications to robustness and approximations for controlled partially observable models; i.e., POMDPs: [43] has studied controlled filter stability where it was shown that one-step observability introduced here leads to stochastic observability universal over all admissible control policies, which then leads to refined robustness results when compared with [32]. In this paper, we consider the control-free case, which allows us to consider N -step observability, with $N > 1$. Additionally, we present numerous explicit examples, which, in the one-step observable setup, is then directly applicable to such robustness results. Recently, filter stability results under total variation (as well as weak convergence under slightly more restrictive setups) have been shown to be consequential in showing the optimality of finite memory control policies in Partially Observed Markov Decision Processes (POMDPs); see [34, Section 4.3 and Theorem 9] and [33, Theorems 3.2, 3.3 and 4.1] where connections with weak merging and total variation merging are made explicit in the approximation error bounds (see [29] for an earlier study where the dependence on filter stability is implicit; further related recent studies include [24]). Accordingly, the results in this paper, notably Theorems II.18 and II.19, are directly applicable in showing that with merging under total variation in expectation, one can show that optimal policies for POMDPs can be approximated by those which use only finite window of measurements and control actions.

III. OBSERVABLE SYSTEM AND MEASUREMENT CHANNEL EXAMPLES

We note that in this section, it will be more convenient to describe our measurement channels via the equivalent functional realization (see (1)), with explicit noise variable Z_n and a measurement function $Y_n = h(X_n, Z_n)$, and thus this will be the convention we will use to define the measurement channel G for the examples presented in the following.

A. Finite state and noise space

Consider a finite setup $\mathcal{X} = \{a_1, \dots, a_n\}$ and $\mathcal{Z} = \{b_1, \dots, b_m\}$. Now, assume $h(x, z)$ has K distinct outputs, where $1 \leq K \leq (n)(m)$ and $\mathcal{Y} = \{c_1, \dots, c_K\}$. We note that for such a setup, there is already a sufficient and necessary condition provided in [56, Theorem V.2]. However, we examine this case to show that our definition is equivalent to the sufficient direction of this theorem, which is the notion of observability presented in [53].

For each x , $h_x(\cdot) := h(x, \cdot)$ can be viewed as a partition of \mathcal{Z} , assigning each $b_i \in \mathcal{Z}$ to an output level $c_j \in \mathcal{Y}$. We can track this by the matrix $H_x(i, j) = 1$ if $h_x(b_i) = c_j$ and zero else. Let Q be the $1 \times m$ vector representing the probability measure of the noise. Let us first consider the one

step observability. Let $g(c_i) = \alpha_i$, with $\alpha = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_K \end{bmatrix}$, and $\int_{\mathcal{Z}} g(h(x, z))Q(dz) =: QH_x\alpha$. Therefore, any function $f(x)$

can be expressed as a $n \times 1$ vector and hence the question reduces to finding a vector α so that $f = QH\alpha$, and the system

is one step observable if and only if the matrix $A \equiv \begin{pmatrix} QH_{a_1} \\ \vdots \\ QH_{a_n} \end{pmatrix}$ is rank n .

Consider then N step observability. We wish to solve equations of the form

$$f(x) = \int_{\mathcal{Y}^N} g(y_{[1,N]}) dP^\mu(y_{[1,N]}|x_1 = x) \quad (8)$$

With knowledge of $Q, h(\cdot, \cdot)$ and T we can directly compute the transition kernel for the joint measure $Y_{[1,n]}|X_1$, however the size of this matrix is n by K^n where $|\mathcal{X}| = n, |\mathcal{Y}| = K$ so complexity grows exponentially. We can deduce a sufficient, but not necessary, condition for n step observability using the marginal conditional measures. Consider that $P^\mu(y_k \in \cdot | X_1 = a_j) = T(a_j | \cdot) T^{k-2} A$, $k \geq 2$ where $T(a_j | \cdot)$ represents the j^{th} row of the transition matrix. Note that these are all $1 \times K$ vectors and represent the marginal measures of $Y_k|X_1$. Consider the class of functions $\mathcal{G}^n = \{g : \mathcal{Y}^n \rightarrow \mathbb{R}\}$ and a subclass $\mathcal{G}_{LC}^n = \{g(y_{[1,n]}) = \sum_{i=1}^n \alpha_i g_i(y_i) | \alpha_i \in \mathbb{R}, g_i \in \mathcal{G}^1\}$. That is, a linear combination of functions of the individual y_i values. We can use these functions to deduce a sufficient, but not necessary, condition for observability.

Lemma III.1. *Assume that $|\mathcal{X}| = n$ and define the matrix*

$$M = (A \quad TA \quad \dots \quad T^{n-1}A)$$

which is $n \times nK$ where $K = |\mathcal{Y}|$. If M is rank n , then the system is n step observable. Furthermore, if M is not rank n , appending more blocks of the form $T^k A$ for $k \geq n$ will not increase the rank of M .

Proof. Begin with (8), consider a restriction to \mathcal{G}_{LC}^n , that is we require g to be of the form $g(y_{[1,n]}) = \sum_{i=1}^n g_i(y_i)$. Denote the $(nK) \times 1$ vector

$\alpha = (g_1(c_1), \dots, g_1(c_K), \dots, g_n(c_1), \dots, g_n(c_K))$. Then

$$\begin{aligned} f(x) &= \sum_{i=1}^n P^\mu(y_i \in \cdot | X_1 = x) \begin{pmatrix} g_i(c_1) \\ \vdots \\ g_i(c_K) \end{pmatrix} \\ &= (QH_x \quad T(x|\cdot)A \quad \dots \quad T(x|\cdot)T^{n-2}A) \alpha \end{aligned}$$

We can then see that this matrix is the j^{th} row of M when $x = a_j$, therefore we have $\begin{pmatrix} f(a_1) \\ \vdots \\ f(a_n) \end{pmatrix} = (A \quad TA \quad \dots \quad T^{n-1}A) \alpha$. If M is rank n , then any function $f : \mathcal{X} \rightarrow \mathbb{R}$ can be expressed as a vector g put through matrix M and the system is observable.

Consider if M is not rank n and if we append another block $T^n A$ to M . By the Cayley-Hamilton theorem, T^n is a linear combination of lower powers of T , e.g. $T^n = \sum_{i=0}^{n-1} \alpha_i T^i$ for some coefficients α_i . Therefore this additional block is a linear combination of the previous blocks, and adds no dimension to the matrix M . \square

If the conditions of this lemma fail, i.e. M is not rank n , that means integrating g over the marginal measures cannot generate any f function. Yet the product of the marginal measures is not the joint measure since the $Y_i|X_1$ are not independent. Hence, working with the marginal measures only is not enough to determine observability as also noted in [53, Remark 13] in a slightly different setup.

Consider the following example. Let $\mathcal{X} = \{1, 2, 3, 4\}$ and $Y = 1_{X \leq 2}$. This can be realized as

$$A = \begin{pmatrix} QH_1 \\ \vdots \\ QH_4 \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 0 \\ 1 & 0 \end{pmatrix}$$

Consider the following transition kernel,

$$T = \begin{pmatrix} 0 & \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \\ \frac{1}{2} & 0 & 0 & \frac{1}{2} \\ 0 & \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \\ \frac{1}{2} & 0 & 0 & \frac{1}{2} \end{pmatrix}$$

Notice that the odd and even rows are identical. If we consider the marginal measures of $Y_1|X_1, \dots, Y_4|X_1$ we have the matrix

$$(A \quad \dots \quad T^3 A) = \begin{pmatrix} 0 & 1 & 0.75 & 0.25 & 0.5625 & 0.4375 & 0.609375 & 0.390625 \\ 0 & 1 & 0.50 & 0.50 & 0.6250 & 0.3750 & 0.593750 & 0.406250 \\ 1 & 0 & 0.75 & 0.25 & 0.5625 & 0.4375 & 0.609375 & 0.390625 \\ 1 & 0 & 0.50 & 0.50 & 0.6250 & 0.3750 & 0.593750 & 0.406250 \end{pmatrix}$$

which is only rank 3, not rank 4. Therefore, we cannot use the marginal measures to determine observability.

However, if we consider the joint measure of $(Y_1, Y_2)|X_1$ we have the matrix

$$A' = \begin{pmatrix} 0 & 0 & \frac{3}{4} & \frac{1}{4} \\ 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{3}{4} & \frac{1}{4} & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 & 0 \end{pmatrix}$$

Where row i is conditioned on $x = i$ and the columns are ordered in binary $y_2 y_1$, e.g. $P(y_1 = 1, y_2 = 0 | x_1 = 2)$ is row 2 column 3. This matrix is full rank, hence the system is N step observable with $N = 2$, even though the marginal measures failed to be full rank.

B. Compact state and noise spaces with affine observations

Consider \mathcal{X}, \mathcal{Z} as compact subsets of \mathbb{R} and let $h(x, z) = a(z)x + b(z)$ for some functions a, b where the image of \mathcal{Z} under a and b is compact (this ensures that \mathcal{Y} is compact). Note that for a fixed choice of z , this is an affine function of x . We will arrive at sufficient conditions for one step observability. Since \mathcal{X} is compact, the set of polynomials is dense in the set of continuous and bounded functions. Therefore, rather than working with a function $f \in C_b(\mathcal{X})$ without loss of generality we assume f is a polynomial. Let $\mathcal{M}_b(\mathbb{R})$ represent the measurable and bounded functions on the real line and consider the mapping

$$S : \mathcal{M}_b(\mathbb{R}) \rightarrow C_b(\mathbb{R}) \quad S(g)(\cdot) \mapsto \int_{\mathcal{Z}} g(h(\cdot, z)) Q(dz)$$

Let $\mathbb{R}[x]_n$ represent the polynomials on the real line up to degree n . Then we have that $S(g)$ is invariant on $\mathbb{R}[x]_n$, that is if g is polynomial of degree n then $S(g)$ is a polynomial of degree n . Furthermore, the coefficients of $S(g)(x) = \sum_{i=0}^n \beta_i x^i$ can be related to the coefficients of $g(x) = \sum_{i=0}^n \alpha_i x^i$ by a linear transformation. Define $N(i, k) = \binom{i}{k} E(a(Z)^k b(Z)^{i-k})$ then by recursive application of binomial theorem we have

$$\begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{pmatrix} = \begin{pmatrix} N(0,0) & N(1,0) & \cdots & N(n,0) \\ 0 & N(1,1) & \cdots & N(n,1) \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & 0 & N(n,n) \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{pmatrix}$$

if we want to generate any polynomial, we require this matrix to be invertible, and since it is upper triangular this amounts to none of the diagonal entries being zero, that is $E[a(z)^n] \neq 0 \forall n \in \mathbb{N}$. Furthermore, we want g to be bounded so we must have $N(n, k) < \infty \forall n \in \mathbb{N}, k \in \{0, \dots, i\}$.

Example

Consider $\mathcal{X} = [-10, 10]$, $\mathcal{Z} = [-1, 1]$, $Z \sim \text{Uni}([-1, 1])$ and $y = z^2 x + z$. We then have $\mathcal{Y} = [-11, 11]$. For any $n \in \mathbb{N}$ we have

$$E[a(z)^n] = \frac{1}{2} \int_{-1}^1 z^{2n} dz = \frac{1}{2n+1} \neq 0$$

additionally, for any $n \in \mathbb{N}, k \in \{0, \dots, n\}$ we have

$$\begin{aligned} N(n, k) &= \binom{n}{k} E(a(z)^k b(z)^{n-k}) = \binom{n}{k} E(z^{n+k}) \\ &= \binom{n}{k} \frac{1}{n+k+1} < \infty \end{aligned}$$

C. A non-linear measurement function

Consider \mathcal{X} as a compact subset of \mathbb{R} , $\mathcal{Z} = \mathbb{R}$. Let $h(x, z) = 1_{x>z}x + 1_{x \leq z}z$ and assume that Q admits a density with respect to Lebesgue. We have

$$\int_{\mathcal{Z}} g(h(x, z))Q(dz) = \int_{-\infty}^x g(x)q(z)dz + \int_x^{\infty} g(z)q(z)dz$$

again, we can approximate any continuous and bounded function f on \mathcal{X} as polynomial, so we assume f is differentiable. We have

$$\begin{aligned} f(x) &= \int_{-\infty}^x g(x)q(z)dz + \int_x^{\infty} g(z)q(z)dz \\ f'(x) &= g(x)q(x) + \int_{-\infty}^x g'(x)q(z)dz - g(x)q(x) \\ &= g'(x)Q(Z \leq x) \end{aligned}$$

Since \mathcal{X} is compact there exists some $x_{\min} \in \mathbb{R}$ such that $x_{\min} < x, \forall x \in \mathcal{X}$. We require for some $\epsilon > 0$ that $Q(Z < x_{\min}) > \epsilon$. This condition says every $x \in \mathcal{X}$ has some positive probability of being observed through $h(x, z)$ and we will not always get pure noise. Then we have

$$g'(x) = 1_{\mathcal{X}}(x) \frac{f'(x)}{Q(Z \leq x)}$$

$$g(x) = c + \int_{-\infty}^x 1_{\mathcal{X}}(u) \frac{f'(u)}{Q(Z \leq u)} du$$

for some constant c . Therefore, we only need to define g over \mathcal{X} . Furthermore, we require g to be bounded, which is implied if g' is bounded since g is only defined over a compact space.

D. Local observability for a non-compact state space

We now study a system which does not have a compact state signal space and satisfies the definitions of local predictability and local observability, so that we can apply Theorem II.15. Consider the POMP with the following transition and measurement kernels

$$\begin{aligned} X_{n+1} &= X_n + N(1, 1) \\ Y_n &= \begin{cases} X_n + 1 & w.p. \frac{1}{2} \\ X_n - 1 & w.p. \frac{1}{2} \end{cases} \end{aligned}$$

We will first show this system is locally predictable. Given an observation Y_{n-1} , it must be that $X_{n-1} = Y_{n-1} - 1$ or $Y_{n-1} + 1$, therefore any filter at time $n-1$ will consist of two point masses at $Y_{n-1} - 1$ and $Y_{n-1} + 1$ with the probability of these two points dependent on the prior. Therefore the predictor at time n will be a convex combination of Gaussian random variables $\alpha_n \mathcal{N}(Y_{n-1}, 1) + (1 - \alpha_n) \mathcal{N}(Y_{n-1} + 2, 1)$ where α_n is determined by the prior.

However, regardless of the value of α for any $\epsilon > 0$ have some compact set K_ϵ such that $\pi_{n-}^\nu(K_\epsilon + Y_{n-1}) > 1 - \epsilon$ for any choice of ν . Therefore the system is locally predictable.

For local observability, assume K is an interval $[-M, M]$ for some whole number $M > 0$ and pick a centering value a . Fix a continuous and bounded function f . We wish to demonstrate a function g that approximates f well over $K+a$ when integrated over the the measurement channel. g must be bounded with a bound that does not depend on a . If we define $g(y)$ recursively as follows:

$$g(y) = \begin{cases} 0 & y < -M + a + 1 \\ 2f(y-1) & y \in [-M + a + 1, -M + a + 3] \\ 2f(y-1) - g(y-2) & y \in [-M + a + 3, M + a + 1] \\ -g(y-2) & y > M + a + 1 \end{cases}$$

g is akin to a telescoping sum in that it cancels out its own previous values. We have

$$\int g(h(x, z))Q(dz) = \frac{1}{2}(g(x+1) + g(x-1))$$

For $x < -M + a$ we have $x-1 < x+1 < -M + a + 1$ hence $g(x-1) = g(x+1) = 0$. For $x \in [-M + a, -M + a + 2]$ we have $g(x+1) = 2f(x+1-1) = 2f(x)$ while $g(x-1) = 0$. Then for $x \in [-M + a + 2, M + a]$ we have

$$g(x+1) = 2f(x+1-1) - g(x+1-2) = 2f(x) - g(x-1)$$

which will cancel with the other $g(x-1)$ term, hence the telescoping. For $x > M + a$ we have $g(x+1) = -g(x+1-2) = -g(x-1)$ hence it will cancel with the previous value.

In each iteration of telescoping, $\|g\|_\infty$ increases by at most $2\|f\|_\infty$, there are $2M$ iterations of telescoping so the overall bound on $\|g\|_\infty$ is $4M\|f\|_\infty$. Therefore we have

$$\|g\|_\infty \leq 4M\|f\|_\infty$$

$$\begin{aligned} \int g(h(x, z)Q(dz)) &= f(x) & x \in [-M + a, M + a] \\ \left| \int g(h(x, z)Q(dz)) \right| &= 0 & x \notin [-M + a, M + a] \end{aligned}$$

this proves local observability.

IV. PROOFS

A. Observability: Proof of Theorem II.7

Lemma IV.1. *Let g be a bounded and measurable function on $(\mathcal{Y}^{k+1}, \mathcal{B}(\mathcal{Y}^{k+1}))$. For any initial prior μ we have*

$$\begin{aligned} & \int_{\mathcal{Y}^{k+1}} g(y_{[n, n+k]}) P^\mu(dy_{[n, n+k]} | Y_{[0, n-1]}) \\ &= \int_{\mathcal{X}} \int_{\mathcal{Y}^{k+1}} g(y_{[n, n+k]}) P(dy_{[n, n+k]} | X_n = x_n) \pi_{n-}^\mu(dx_n) \end{aligned} \quad (9)$$

Proof.

$$\begin{aligned} & \int_{\mathcal{Y}^{k+1}} g(y_{[n, n+k]}) P^\mu(dy_{[n, n+k]} | Y_{[0, n-1]}) \\ &= \int_{\mathcal{Y}^{k+1} \times \mathcal{X}} g(y_{[n, n+k]}) P^\mu(d(y_{[n, n+k]}, x_n) | Y_{[0, n-1]}) \end{aligned}$$

we then apply the chain rule for conditional probability measures and we have

$$\int_{\mathcal{X}} \int_{\mathcal{Y}^{k+1}} g(y_{[n, n+k]}) P^\mu(dy_{[n, n+k]} | x_n, Y_{[0, n-1]}) \pi_{n-}^\mu(dx_n)$$

Since $\{(X_n, Y_n)\}_{n=0}^\infty$ is a Markov chain chain, $Y_{[n, n+k]}$ is conditionally independent of $Y_{[0, n-1]}$ given X_n . Additionally, the prior does not determine the conditional measure, therefore we have

$$\int_{\mathcal{X}} \int_{\mathcal{Y}^{k+1}} g(y_{[n, n+k]}) P(dy_{[n, n+k]} | x_n) \pi_{n-}^\mu(dx_n)$$

where we do not include a prior in the superscript of the conditional measure, since the conditional measure is the same regardless of the prior. \square

Corollary IV.2. *Let g be a bounded and measurable function on $(\mathcal{Y}, \mathcal{B}(\mathcal{Y}))$. For any prior μ we have*

$$\int_{\mathcal{Y}} g(y_n) P^\mu(dy_n | X_n = x) = \int_{\mathcal{Z}} g(h_x(z)) Q(dz) \quad (10)$$

Proof. Z is a random variable on the probability space $(\mathcal{Z}, \mathcal{B}(\mathcal{Z}), Q)$ and Y_n exists on the measurable space $(\mathcal{Y}, \mathcal{B}(\mathcal{Y}))$. Then, for every fixed choice of $X_n = x$ we have that Y_n is a fixed function of Z , that is $Y_n = h_x(Z)$. For any set $A \in \mathcal{B}(\mathcal{Y})$ we have $P^\mu(Y_n \in A | X_n = x) = Q(h_x^{-1}(A))$. Yet this means that $P^\mu(Y_n \in \cdot | X_n = x)$ is exactly the pushforward measure of Q under the mapping h_x , call this measure $h_x Q(A) = Q(h_x^{-1}(A))$. We then have:

$$\int_{\mathcal{Y}} g(y) h_x Q(dy) = \int_{\mathcal{Z}} g(h_x(z)) Q(dz). \quad \square$$

Notice that the inner integral in the RHS of Equation (9) is a function of x . The LHS is then the term considered

in the total variation merging of the predictive measures of the measurement sequences, while the RHS is the term considered in the weak merging of the one-step predictor. We can then leverage Blackwell and Dubin's theorem to arrive at a sufficient condition for weak merging of the one-step predictor. Theorem II.7 is closely related to [55, Prop. 3.11] and its proof is in essence a sufficient condition for uniform observability (of the predictor).

Proof of Theorem II.7.

Fix any $f \in C_b(\mathcal{X})$ and $\epsilon > 0$. We wish to show that $\exists N$ such that $\forall n > N$,

$$\left| \int f d\pi_{n-}^\mu - \int f d\pi_{n-}^\nu \right| < \epsilon$$

By observability for the fixed f , (5) holds for some $N' + 1$. Therefore we can find some g with $\|g\|_\infty < \infty$ such that

$$\tilde{f}(x) = \int_{\mathcal{Y}^{N'+1}} g(y_{[1, 1+N']}) P(dy_{[1, 1+N']} | X_1 = x)$$

and $\|f - \tilde{f}\|_\infty < \frac{\epsilon}{3}$. Conditioned on the value of $X_n = x$ and since the noise is i.i.d, the conditional channel $Y_{[n, n+N']} | X_n$ is time invariant, so it holds that

$$\tilde{f}(x) = \int_{\mathcal{Y}^{N'+1}} g(y_{[n, n+N']}) P(dy_{[n, n+N']} | X_n = x)$$

is the same regardless of the choice of n . Then we have

$$\begin{aligned} & \left| \int f d\pi_{n-}^\mu - \int f d\pi_{n-}^\nu \right| \\ & \leq \left| \int \tilde{f} d\pi_{n-}^\mu - \int \tilde{f} d\pi_{n-}^\nu \right| + \left| \int (f - \tilde{f}) d\pi_{n-}^\mu \right| \\ & \quad + \left| \int (f - \tilde{f}) d\pi_{n-}^\nu \right| \end{aligned} \quad (11)$$

Now, by assumption $\|f - \tilde{f}\|_\infty < \frac{\epsilon}{3}$ therefore the last two terms are less than $\frac{2}{3}\epsilon$. We then apply Lemma IV.1 and we have

$$\begin{aligned} & \left| \int \tilde{f} d\pi_{n-}^\mu - \int \tilde{f} d\pi_{n-}^\nu \right| + \frac{2}{3}\epsilon \\ &= \left| \int_{\mathcal{Y}^{N'+1}} g(y_{[n, n+N']}) P^\mu(dy_{[n, n+N']} | Y_{[0, n-1]}) \right. \\ & \quad \left. - \int_{\mathcal{Y}^{N'+1}} g(y_{[n, n+N']}) P^\nu(dy_{[n, n+N']} | Y_{[0, n-1]}) \right| + \frac{2}{3}\epsilon \end{aligned}$$

By Assumption 6, we have $P^\mu(Y_{[0, \infty)} \in \cdot) \ll P^\nu(Y_{[0, \infty)} \in \cdot)$. Then via a classic result by Blackwell and Dubins [7], we have that $P^\mu(Y_{[n, n+N']} \in \cdot | Y_{[0, n-1]})$ and $P^\nu(Y_{[n, n+N']} \in \cdot | Y_{[0, n-1]})$ merge in total variation P^μ a.s. as $n \rightarrow \infty$. Define $\tilde{g} = \frac{g}{\|g\|_\infty}$. Then $\exists N \in \mathbb{N}$ such that $\forall n > N$,

$$\begin{aligned} & \left| \int_{\mathcal{Y}^{N'+1}} \tilde{g}(y_{[n, n+N']}) P^\mu(dy_{[n, n+N']} | Y_{[0, n-1]}) \right. \\ & \quad \left. - \int_{\mathcal{Y}^{N'+1}} \tilde{g}(y_{[n, n+N']}) P^\nu(dy_{[n, n+N']} | Y_{[0, n-1]}) \right| < \frac{\epsilon}{3\|g\|_\infty} \end{aligned}$$

we then have:

$$\left| \int_{\mathcal{Y}^{N'+1}} g(y_{[n, n+N']}) P^\mu(dy_{[n, n+N']} | Y_{[0, n-1]}) \right|$$

$$\begin{aligned} & - \int_{\mathcal{Y}^{N'+1}} g(y_{[n,n+N']}) P^\nu(dy_{[n,n+N']} | Y_{[0,n-1]}) + \frac{2}{3}\epsilon \\ & \leq \|g\|_\infty \frac{\epsilon}{3\|g\|_\infty} + \frac{2}{3}\epsilon = \epsilon \end{aligned}$$

therefore, since f and ϵ are arbitrary we have for any $f \in C_b(S)$: $\lim_{n \rightarrow \infty} |\int f d\pi_{n-}^\mu - \int f d\pi_{n-}^\nu| = 0$, which means π_{n-}^μ and π_{n-}^ν merge weakly. \square

B. Weak Filter Stability: Proof of Theorem II.9

Here we will utilize results from [31]. This paper was concerned with a different topic than filter stability, namely the weak Feller property of the ‘‘filter update’’ kernel. That is, one can view the filter π_n^μ and the measurement Y_n as its own Markov chain $\{(\pi_n^\mu, Y_n)\}_{n=0}^\infty$ which takes values in $\mathcal{P}(\mathcal{X}) \times \mathcal{Y}$. The filter update kernel is the transition kernel of this Markov chain. We will not study this kernel, but some of the analysis in [31] is useful in providing concise arguments to connect the filter to the predictor.

Proof of Theorem II.9.

Begin by assuming that the predictor merges weakly almost surely. As is argued in [31], one can view the filter π_n^μ as a function of π_{n-1}^μ (the previous filter) and the current observation $Y_n = y_n$, that is $\pi_n^\mu = F(\pi_{n-1}^\mu, y_n)$. Pick any continuous and bounded function f , we have

$$\begin{aligned} & E^\mu \left[\left| \int_{\mathcal{X}} f(x) \pi_n^\mu(dx) - \int_{\mathcal{X}} f(x) \pi_n^\nu(dx) \right| \right] \\ & = E^\mu \left[E^\mu \left[\left| \int_{\mathcal{X}} f(x) F(\pi_{n-1}^\mu, y_n)(dx) \right. \right. \right. \\ & \quad \left. \left. \left. - \int_{\mathcal{X}} f(x) F(\pi_{n-1}^\nu, y_n)(dx) \right| \middle| Y_{[0,n-1]} \right] \right] \quad (12) \end{aligned}$$

Now, define the set $I^+(y_{[0,n-1]}) \subset \mathcal{Y}$ as:

$$\begin{aligned} I^+(y_{[0,n-1]}) & = \{y_n \in \mathbb{Y} \mid \int_{\mathcal{X}} f(x) F(\pi_{n-1}^\mu, y_n)(dx) \\ & > \int_{\mathcal{X}} f(x) F(\pi_{n-1}^\nu, y_n)(dx)\} \end{aligned}$$

where the argument $y_{[0,n-1]}$ is the sequence on which the previous filters π_{n-1}^μ and π_{n-1}^ν are realized. Define the complement of this set as $I^-(y_{[0,n-1]})$. Then for every fixed realization $y_{[0,n-1]}$ we can break the inner expectation in (12) (which is an integral) into two parts and follow the analysis in [31, Equation 4] together with Theorem 8.6.2 in [8] to arrive at the conclusion. \square

C. Local Observability: Proof of Theorem II.12 and II.15

The idea of local observability is the shift some of the burden of approximating the signal f . When we work with a function

$$\tilde{f}(x) = \int_{\mathcal{Y}^{N'+1}} g(y_{[n,n+N']}) P(dy_{[n,n+N']} | X_n = x)$$

the result is the terms seen in equation (11). The first term is dealt with by Blackwell and Dubin’s theorem, so we must

make sure the second and third term can be made arbitrarily small. For any set K we can write

$$\begin{aligned} & \left| \int (f - \tilde{f}) d\pi_{n-}^\nu \right| \leq \sup_{x \in K} |f(x) - \tilde{f}(x)| \pi_{n-}^\nu(K) \\ & + \sup_{x \notin K} |f(x) - \tilde{f}(x)| \pi_{n-}^\nu(K^C) \end{aligned}$$

in the previous result we bounded this by simply approximating f well over the whole space. Instead, we can choose a K where \tilde{f} approximates f well over K and $\pi_{n-}^\nu(K^C)$ makes the other term arbitrarily small. Furthermore, by taking advantage of the full supremum of total variation we can work with a series of uniformly bounded functions \tilde{f}_n and shifting sets K_n that change with n .

Proof of Theorem II.12. Pick any continuous and bounded function f and any $\epsilon > 0$. Fix any sequence of observations $y_{[0,\infty)}$ where the predictors π_{n-}^μ and π_{n-}^ν are well defined and maintain this sequence for the remainder of the proof. Then consider

$$\lim_{n \rightarrow \infty} \left| \int_{\mathcal{X}} f(x) \pi_{n-}^\mu(dx) - \int_{\mathcal{X}} f(x) \pi_{n-}^\nu(dx) \right|$$

For any function series of functions \tilde{f}_n of x we have an upper bound

$$\begin{aligned} & \lim_{n \rightarrow \infty} \left| \int_{\mathcal{X}} \tilde{f}_n(x) \pi_{n-}^\mu(dx) - \int_{\mathcal{X}} \tilde{f}_n(x) \pi_{n-}^\nu(dx) \right| \\ & + \left| \int_{\mathcal{X}} (f - \tilde{f}_n)(x) \pi_{n-}^\mu(dx) \right| + \left| \int_{\mathcal{X}} (f - \tilde{f}_n)(x) \pi_{n-}^\nu(dx) \right| \end{aligned}$$

By assumption of K local predictability, we have a compact sets $K_n = K + a_n$ where $\pi_{n-}^\mu(K_n) = 1$ for every $\mu \ll \nu$ and every n .

By K local observability, we can find a uniformly bounded series of functions $g_n \leq M$ where

$$\begin{aligned} \tilde{f}_n(x) & = \int_{\mathcal{Z}} g_n(h(x, z)) Q(dz) \\ \sup_{x \in K_n} |f(x) - \tilde{f}_n(x)| & \leq \frac{\epsilon}{3} \end{aligned}$$

then for the two approximation terms we have

$$\begin{aligned} & \left| \int_{\mathcal{X}} (f - \tilde{f}_n)(x) \pi_{n-}^\nu(dx) \right| \\ & \leq \sup_{x \in K_n} |f(x) - \tilde{f}_n(x)| \pi_{n-}^\nu(K_n) + \sup_{x \notin K_n} |f(x) - \tilde{f}_n(x)| \pi_{n-}^\nu(K_n^C) \\ & \leq \frac{\epsilon}{3} \end{aligned}$$

we then have

$$\begin{aligned} & \lim_{n \rightarrow \infty} \left| \int_{\mathcal{X}} \tilde{f}_n(x) \pi_{n-}^\mu(dx) - \int_{\mathcal{X}} \tilde{f}_n(x) \pi_{n-}^\nu(dx) \right| + \frac{2}{3}\epsilon \\ & = \lim_{n \rightarrow \infty} \left| \int_{\mathcal{Y}} g_n(y_n) P^\mu(dy_n | y_{[0,n-1]}) \right. \\ & \quad \left. - \int_{\mathcal{Y}} g_n(y_n) P^\nu(dy_n | y_{[0,n-1]}) \right| + \frac{2}{3}\epsilon \end{aligned}$$

we must appeal to the full uniform bound of the Blackwell and Dubins theorem, which was not required in the proof of

Theorem II.7. The full statement of the Blackwell and Dubins theorem tells us that

$$\lim_{n \rightarrow \infty} \sup_{\|g\| \leq 1} \left| \int_{\mathcal{Y}} g(y_n) P^\mu(dy_n | y_{[0, n-1]}) - \int_{\mathcal{Y}} g(y_n) P^\nu(dy_n | y_{[0, n-1]}) \right| = 0, \quad (13)$$

where the supremum is taken over measurable functions g . Thus, for any fixed measurable and bounded function g , we have that

$$\left| \int_{\mathcal{Y}} g(y_n) P^\mu(dy_n | y_{[0, n-1]}) - \int_{\mathcal{Y}} g(y_n) P^\nu(dy_n | y_{[0, n-1]}) \right|$$

converges to 0 as $n \rightarrow \infty$; this was the form of the statement utilized in the proof of Theorem II.7. However, if we have a sequence of measurable functions g_n with a uniform bound, $g_n \leq M \forall n \in \mathbb{N}$, then the supremum in (13) allows us to make a uniform claim about the convergence to zero of the sequence,

$$\left| \int_{\mathcal{Y}} g_n(y_n) P^\mu(dy_n | y_{[0, n-1]}) - \int_{\mathcal{Y}} g_n(y_n) P^\nu(dy_n | y_{[0, n-1]}) \right|,$$

and this completes the proof. \square

Proof of Theorem II.15. Fix any f and any ϵ . We begin from the upper bound used previously

$$\lim_{n \rightarrow \infty} \left| \int_{\mathcal{X}} \tilde{f}_n(x) \pi_{n-}^\mu(dx) - \int_{\mathcal{X}} \tilde{f}_n(x) \pi_{n-}^\nu(dx) \right| + \left| \int_{\mathcal{X}} (f - \tilde{f}_n)(x) \pi_{n-}^\mu(dx) \right| + \left| \int_{\mathcal{X}} (f - \tilde{f}_n)(x) \pi_{n-}^\nu(dx) \right|$$

for some series of functions \tilde{f}_n .

By local predictability, the shifted predictors are a tight family. Therefore for any ϵ' we have a series of compact sets $K_n = K' + a_n$ such that $\pi_{n-}^\nu(K_n) \geq 1 - \epsilon'$ for any $\mu \ll \nu$ and any n .

The proof then proceeds similarly as that of Theorem II.12 \square

D. Predictor Merging in Total Variation: Proof of Theorem II.18

We now extend our results from weak merging to total variation. We first state the following supporting results.

Lemma IV.3. *The (measurement-update) map:*

$$(\pi_{n-}, y) \mapsto \pi_n \quad : \quad \pi_n(\cdot) := E_{\pi_{n-}}[1_{X_n \in \cdot} | Y_n = y]$$

which maps from $\mathcal{P}(\mathcal{X}) \times \mathbb{Y}$ to $\mathcal{P}(\mathcal{X})$ is weakly continuous in π_{n-} for almost every y , provided that $g(x, y)$ is positive, bounded and continuous in x for every fixed y .

Proof. Consider a continuous and bounded f and let $\pi_{n-}^m \rightarrow \pi_{n-}$ weakly. Then,

$$\begin{aligned} E_{\pi_{n-}^m}[f(x_n) | Y_n = y_n] &= \int f(x_n) \frac{g(x_n, y_n) \pi_{n-}^m(dx_n)}{\int_{\mathcal{X}} g(x_n, y_n) \pi_{n-}^m(dx_n)} \\ &= \frac{\int f(x_n) g(x_n, y_n) \pi_{n-}^m(dx_n)}{\int_{\mathcal{X}} g(x_n, y_n) \pi_{n-}^m(dx_n)} \end{aligned}$$

Since $g(\cdot, y_n)$ is bounded and continuous, both the numerator and the denominator converge. \square

Lemma IV.4. *Let $T(dx_1 | x) = t(x_1, x) \phi(dx_1)$ where t is continuous in x for every x_1 . Then, the (time-update) map:*

$$(\pi_n) \mapsto \pi_{n+1-} \quad : \quad \pi_{n+1-}(\cdot) := \int T(\cdot | x_n) \pi_n(dx_n)$$

which maps from $\mathcal{P}(\mathcal{X})$ to $\mathcal{P}(\mathcal{X})$ is so that if $\pi_n^\nu \rightarrow \pi_n^\mu$ weakly then $\pi_{n+1-}^\nu \rightarrow \pi_{n+1-}^\mu$ in total variation.

Proof. We will build on Scheffé's Lemma [5]. For every given history, we have

$$\pi_{n+1-}^\nu(dx_{n+1}) = \int T(dx_{n+1} | x_n) \pi_n^\nu(dx_n)$$

Now, $\int T(dx_{n+1} | x_n)$ is so that,

$$\int t(x_{n+1}, x_n) \phi(dx_{n+1}) \pi_n^m(dx_n) \rightarrow \int t(x_{n+1}, x_n) \phi(dx_{n+1}) \pi_n(dx_n)$$

in total variation since for every fixed z , the Radon-Nikodym derivative (density) with respect to ϕ

$$\frac{\int t(x_{n+1}, x_n) \phi(\cdot) \pi_n^m(dx_n)}{d\phi}(z) = \int t(z, x_n) \pi_n^m(dx_n)$$

satisfies pointwise convergence

$$\int t(z, x_n) \pi_n^\nu(dx_n) \rightarrow \int t(z, x_n) \pi_n^\mu(dx_n)$$

and Scheffé's lemma implies that convergence is in total variation. Now, we can apply the result to the sequence π_n^ν converging to π_n^μ . \square

Proof. Proof of Theorem II.18(i)

Under Assumption II.16, the proof follows from Lemma IV.3 and IV.4. While in Lemma IV.3 and IV.4 we consider convergence (and not merging), we note that the proof of Lemma IV.3 also implies weak merging of the posteriors as the priors weakly merge, and by considering the signed measure $\pi_n^{\nu, \gamma} - \pi_n^{\mu, \gamma}$ in the proof of Lemma IV.4, total variation merging is a result of a generalized Scheffé's lemma [8, Theorem 2.8.9]. \square

Lemma IV.5. *Let \exists some measure $\bar{\mu}$ such that $T(\cdot | x) \ll \bar{\mu}$ for every $x \in \mathcal{X}$. Then we have that $\pi_{n-}^\mu, \pi_{n-}^\nu \ll \bar{\mu}$ for every $n \in \mathbb{N}$*

Proof. For all $n \geq 1$ we have

$$\begin{aligned} \pi_{n-}^\mu(A) &= \int_{\mathcal{X}} T(A | x) \pi_{n-1}^\mu(dx) = \int_{\mathcal{X}} \int_A \frac{dT(\cdot | x)}{d\bar{\mu}}(a) \bar{\mu}(da) \pi_{n-1}^\mu(dx) \\ &= \int_A \left(\int_{\mathcal{X}} \frac{dT(\cdot | x)}{d\bar{\mu}}(a) \pi_{n-1}^\mu(dx) \right) \bar{\mu}(da) \end{aligned}$$

where we have applied Fubini's theorem in the final equality. Therefore π_{n-}^μ is absolutely continuous with respect to $\bar{\mu}$ for every $n \geq 1$. \square

Lemma IV.6. *Let Assumption II.17 hold and let f_{n-}^μ denote the density function of π_{n-}^μ . Fix any sequence of measurements $y_{[0, \infty)}$ and denote the collection of probability density functions $\mathcal{F}^\mu = \{f_{n-}^\mu | n \in \mathbb{N}\}$, $\mathcal{F}^\nu = \{f_{n-}^\nu | n \in \mathbb{N}\}$. Then $\mathcal{F}^\mu, \mathcal{F}^\nu$ are uniformly bounded equicontinuous families.*

Proof. As we see from Lemma IV.5,

$$f_{n-}^{\mu}(x_n) = \frac{d\pi_{n-}^{\mu}}{d\phi}(x_n) = \int_{\mathcal{X}} t(x_n|x_{n-1})\pi_{n-1}^{\mu}(dx_{n-1})$$

Where $t(\cdot|x)$ is the Radon Nikodym derivative of $T(\cdot|x)$ with respect to our dominating measure ϕ and $d(\cdot, \cdot)$ will represent the metric on \mathcal{X} (recall \mathcal{X} is a complete, separable, metric space). We require $\forall \epsilon > 0$, $x^* \in \mathcal{X} \exists \delta > 0$ such that $\forall d(x, x^*) < \delta$, $\forall n \in \mathbb{N}$ we have $|f_{n-}^{\mu}(x) - f_{n-}^{\mu}(x^*)| < \epsilon$. By Assumption II.17, clearly f_{n-}^{μ} is uniformly bounded since t is uniformly bounded. Then, for any $\epsilon > 0$, $\forall x^* \in \mathcal{X}$ we can find a $\delta > 0$ such that $\forall x_1 \in \mathcal{X}$, $|t(x_2|x_1) - t(x^*|x_1)| < \epsilon$ when $d(x_2, x^*) < \delta$. Now, assume $d(x_2, x^*) < \delta$, we have

$$\begin{aligned} |f_{n-}^{\mu}(x_2) - f_{n-}^{\mu}(x^*)| &= \left| \int_{\mathcal{X}} t(x_2|x_1) - t(x^*|x_1) d\pi_{n-}^{\mu}(dx_1) \right| \\ &\leq \int_{\mathcal{X}} |t(x_2|x_1) - t(x^*|x_1)| d\pi_{n-}^{\mu}(x_1) \leq \epsilon \end{aligned}$$

which proves that \mathcal{F}^{μ} and \mathcal{F}^{ν} are uniformly bounded and equicontinuous families. \square

Proof of Theorem II.18(ii). By assumption we have weak stability of the predictor P^{μ} a.s.. Then there exists a set of measure sequences $B \subset \mathcal{Y}_{[0,\infty]}^{\mathbb{Z}^+}$ with $P^{\mu}(B) = 1$. For each measurement sequence $y_{[0,\infty]} \in B$, we have that the predictor realizations π_{n-}^{μ} and π_{n-}^{ν} merge in the weak sense. We will choose a general measurement sequence $y_{[0,\infty]} \in B$ and fix this sequence for the remainder of the proof. Via Lemma IV.5, and IV.6, \mathcal{F}^{μ} and \mathcal{F}^{ν} are uniformly bounded and equicontinuous families. Let $\mathcal{F}^{\mu-\nu} = \{f_n | f_n = f_{n-}^{\mu} - f_{n-}^{\nu}\}$, then the sequence $\{f_n\}_{n=1}^{\infty}$ is a uniformly bounded and equicontinuous class of integrable functions. As in the proof of [40, Lemma 2], now pick a sequence of compact sets $K_j \subset \mathcal{X}$ such that $K_j \subset K_{j+1}$. By the Arzela-Ascoli theorem [48], for any subsequence we can find further subsequences $f_{n_k^j}$ such that

$$\lim_{k \rightarrow \infty} \sup_{x \in K_j} |f_{n_k^j}(x) - f^j(x)| = 0$$

for some continuous function $f^j : K_j \rightarrow [0, \infty)$. Via the K_j being nested, we can have $\{f_{n_k^{j+1}}\}$ be a subsequence of $\{f_{n_k^j}\}$, and therefore $f^{j+1} = f^j$ over K_j . Then define the function \tilde{f} on \mathcal{X} by $\tilde{f}(x) = f^j(x)$, $x \in K_j$. Using Cantor's diagonal method, we can find an increasing sequence of integers $\{m_i\}$ which is a subsequence of $\{n_k^j\}$ for every j . Therefore

$$\lim_{i \rightarrow \infty} f_{m_i}(x) = \tilde{f}(x) \quad \forall x \in \mathcal{X}$$

and the convergence is uniform over each K_j and \tilde{f} is continuous. Now, f_{m_i} converges weakly to the zero measure by assumption, and via uniform convergence for any Borel set B we have

$$\int_B f_{m_i}(x) dx \rightarrow \int_B \tilde{f}(x) dx,$$

i.e. setwise convergence. Yet this implies weak convergence, so $\tilde{f} = 0$ almost everywhere, yet \tilde{f} is continuous so it is 0 everywhere. Now, via Prokhorov theorem (Theorem 8.6.2 in

[8]) we have that $\mathcal{F}^{\mu-\nu}$ is a tight family. Therefore, for every $\epsilon > 0$ we can find a compact set K_{ϵ} such that

$$|\pi_{n-}^{\mu} - \pi_{n-}^{\nu}|(\mathcal{X} \setminus K_{\epsilon}) < \epsilon \quad \forall n \in \mathbb{N}.$$

then we have

$$\begin{aligned} \lim_{i \rightarrow \infty} \|\pi_{m_i-}^{\mu} - \pi_{m_i-}^{\nu}\|_{TV} &\leq \lim_{i \rightarrow \infty} |\pi_{m_i-}^{\mu} - \pi_{m_i-}^{\nu}|(\mathcal{X} \setminus K_{\epsilon}) \\ &\quad + |\pi_{m_i-}^{\mu} - \pi_{m_i-}^{\nu}|(K_{\epsilon}) \\ &\leq \lim_{i \rightarrow \infty} \sup_{\|g\|_{\infty} \leq 1} \left| \int_{K_{\epsilon}} g(x) f_{m_i}(x) dx \right| + \epsilon \\ &\leq \lim_{i \rightarrow \infty} \sup_{\|g\|_{\infty} \leq 1} \left| \int_{K_{\epsilon}} g(x) (\tilde{f} - f_{m_i})(x) dx \right| + \left| \int_{K_{\epsilon}} g(x) \tilde{f}(x) dx \right| + \epsilon \\ &\leq \lim_{i \rightarrow \infty} \|\tilde{f} - f_{m_i}\|_{\infty} \phi(K_{\epsilon}) + \epsilon \end{aligned}$$

since we have already argued $\tilde{f} = 0$. Now, over the compact set K_{ϵ} , f_{m_i} converges to \tilde{f} uniformly, therefore $\exists N$ such that $\forall k > N$, $\|\tilde{f} - f_{n_k}\|_{\infty} < \frac{\epsilon}{\phi(K_{\epsilon})}$. We then conclude that

$$\lim_{i \rightarrow \infty} \|\pi_{m_i-}^{\mu} - \pi_{m_i-}^{\nu}\|_{TV} = 0$$

Thus, for every subsequence of $\{f_n\}_{n=1}^{\infty}$, we can find a subsequence that converges in total variation, which implies that the original sequence converges in total variation. \square

E. Filter Merging in Total Variation: Proof of Theorem II.19

For completeness, in Section VI some supporting results are presented.

Proof of Theorem II.19. The sigma fields $\mathcal{F}_{n,\infty}^{\mathcal{X}} \vee \mathcal{F}_{0,\infty}^{\mathcal{Y}}$ are a decreasing sequence, that is $\mathcal{F}_{n+1,\infty}^{\mathcal{X}} \vee \mathcal{F}_{0,\infty}^{\mathcal{Y}} \subset \mathcal{F}_{n,\infty}^{\mathcal{X}} \vee \mathcal{F}_{0,\infty}^{\mathcal{Y}}$. Therefore, when we take their intersection, removing the first or largest sigma field $\mathcal{F}_{0,\infty}^{\mathcal{X}} \vee \mathcal{F}_{0,\infty}^{\mathcal{Y}}$ from the intersection of a decreasing set of sigma fields does not change the overall intersection. From Lemma VI.5 and VI.8, it is clear that the two conditions for merging in total variation in expectation are equivalent since the sigma fields on the LHS of Equation (16) and (19) are equal. \square

We have now established that the filter merges in total variation in expectation, but we would like to extend this result to almost surely. By a simple application of Fatou's lemma, we can argue the liminf of the total variation of the filter is zero P^{μ} a.s.. Hence if the limit exists, it must be zero, yet it is not immediate that the limit will exist. This leads to the following.

Theorem IV.7. [52, p. 572] *Assume the filter is stable in total variation in expectation. Then the filter is stable in total variation P^{μ} a.s.*

F. Relative Entropy Merging: Proof of Theorem II.20

We will now show that the relative entropy merging of the filter is essentially equivalent to merging in total variation in expectation. Via Lemma VI.4 and VI.6, it is clear that the filter and predictor admit Radon-Nikodym derivatives. Therefore, working with $D(\pi_n^{\mu} || \pi_n^{\nu})$ and $D(\pi_{n-}^{\mu} || \pi_{n-}^{\nu})$ is well defined. A well known result for relative entropy is the chain rule [25, Theorem 5.3.1]:

Lemma IV.8. For joint measures P, Q on random variables X, Y we have

$$D(P(X, Y) \| Q(X, Y)) = D(P(X) \| Q(X)) + D(P(Y|X) \| Q(Y|X))$$

Note for two sigma fields \mathcal{F} and \mathcal{G} and two joint measures P and Q on $\mathcal{F} \vee \mathcal{G}$ one could also express this relationship as

$$D(P|_{\mathcal{F} \vee \mathcal{G}} \| Q|_{\mathcal{F} \vee \mathcal{G}}) = D(P|_{\mathcal{F}} \| Q|_{\mathcal{F}}) + D(P|_{\mathcal{G}} \| Q|_{\mathcal{G}}) \quad (14)$$

Proof of Theorem II.20. First assume the filter is stable in relative entropy. Since the square root function is continuous and convex, we have

$$\begin{aligned} 0 &= \lim_{n \rightarrow \infty} \sqrt{\frac{2}{\log(e)} E^\mu [D(\pi_n^\mu \| \pi_n^\nu)]} \\ &\geq \lim_{n \rightarrow \infty} E^\mu \left[\sqrt{\frac{2}{\log(e)} D(\pi_n^\mu \| \pi_n^\nu)} \right] \end{aligned}$$

where we have applied Jensen's inequality. We then apply Pinsker's inequality and we have $\lim_{n \rightarrow \infty} E^\mu [\|\pi_n^\mu - \pi_n^\nu\|_{TV}] = 0$.

For the converse direction, by chain rule (14), it is clear that

$$\begin{aligned} E^\mu [D(\pi_n^\mu \| \pi_n^\nu)] &= D(P^\mu|_{\mathcal{F}_n^X \vee \mathcal{F}_{0,n}^Y} \| P^\nu|_{\mathcal{F}_n^X \vee \mathcal{F}_{0,n}^Y}) \\ &= D(P^\mu|_{\mathcal{F}_n^X \vee \mathcal{F}_{0,n}^Y} \| P^\nu|_{\mathcal{F}_n^X \vee \mathcal{F}_{0,n}^Y}) - D(P^\mu|_{\mathcal{F}_{0,n}^Y} \| P^\nu|_{\mathcal{F}_{0,n}^Y}) \end{aligned}$$

by the Markov Property we have $X_{[0, n-1]}, Y_{[0, n-1]}$ and $X_{[n+1, \infty)}, Y_{[n+1, \infty)}$ are conditionally independent given X_n, Y_n therefore we have:

$$\begin{aligned} &D(P^\mu|_{\mathcal{F}_n^X \vee \mathcal{F}_{0,n}^Y} \| P^\nu|_{\mathcal{F}_n^X \vee \mathcal{F}_{0,n}^Y}) \\ &= D(P^\mu|_{\mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y} \| P^\nu|_{\mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y}) \end{aligned}$$

Then $\mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y$ is a decreasing sequence of sigma fields. By [4, Theorem 2] we have that if the relative entropy is ever finite, the limit of the relative entropy restricted to these sigma fields is the relative entropy restricted to the intersection of the decreasing fields, that is

$$\begin{aligned} &\lim_{n \rightarrow \infty} D(P^\mu|_{\mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y} \| P^\nu|_{\mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y}) \\ &= D(P^\mu|_{\bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y} \| P^\nu|_{\bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y}) \end{aligned}$$

Likewise, $\mathcal{F}_{0,n}^Y$ is an increasing sequence of sigma fields, therefore by [4, Theorem 3] we have that if the relative entropy is ever finite, the relative entropy restricted to these sigma fields is the relative entropy over the limit field, that is

$$\lim_{n \rightarrow \infty} D(P^\mu|_{\mathcal{F}_{0,n}^Y} \| P^\nu|_{\mathcal{F}_{0,n}^Y}) = D(P^\mu|_{\mathcal{F}_{0,\infty}^Y} \| P^\nu|_{\mathcal{F}_{0,\infty}^Y})$$

Therefore,

$$\begin{aligned} &\lim_{n \rightarrow \infty} E^\mu [D(\pi_n^\mu \| \pi_n^\nu)] \\ &= D(P^\mu|_{\bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y} \| P^\nu|_{\bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y}) \\ &\quad - D(P^\mu|_{\mathcal{F}_{0,\infty}^Y} \| P^\nu|_{\mathcal{F}_{0,\infty}^Y}) \end{aligned}$$

By Lemma VI.1 we have

$$\frac{dP^\mu|_{\bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y}}{dP^\nu|_{\bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y}} = E^\nu \left[\frac{d\mu}{d\nu}(X_0) \Big| \bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y \right] = f_1$$

$$\frac{dP^\mu|_{\mathcal{F}_{0,\infty}^Y}}{dP^\nu|_{\mathcal{F}_{0,\infty}^Y}} = E^\nu \left[\frac{d\mu}{d\nu}(X_0) \Big| \mathcal{F}_{0,\infty}^Y \right] = f_2$$

Note that f_1 is $\bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y$ measurable, while f_2 is $\mathcal{F}_{0,\infty}^Y$ measurable, and $\mathcal{F}_{0,\infty}^Y \subset \bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y$. By Lemma VI.5, we have that if the filter merges in total variation in expectation, then for a set of state and observation sequences $\omega = (x_i, y_i)_{i=0}^\infty \in A \subset \mathcal{F}_{0,\infty}^X \vee \mathcal{F}_{0,\infty}^Y$ with $P^\nu(A) = 1$, we have $f_1(\omega) = f_2(\omega)$. Yet this then means over the set A of P^ν measure 1, $f_1 = f_2$ is $\mathcal{F}_{0,\infty}^Y$ measurable. We then have

$$\begin{aligned} &D(P^\mu|_{\bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y} \| P^\nu|_{\bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y}) \\ &\quad - D(P^\mu|_{\mathcal{F}_{0,\infty}^Y} \| P^\nu|_{\mathcal{F}_{0,\infty}^Y}) \\ &= E^\mu [\log(f_1)] - E^\mu [\log(f_2)] = E^\nu [f_1 \log(f_1)] - E^\nu [f_2 \log(f_2)] \\ &= \int_\Omega f_1(\omega) \log(f_1(\omega)) dP^\nu|_{\bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y}(\omega) \\ &\quad - \int_\Omega f_2(\omega) \log(f_2(\omega)) dP^\nu|_{\mathcal{F}_{0,\infty}^Y}(\omega) \\ &= \int_A f_1(\omega) \log(f_1(\omega)) dP^\nu|_{\bigcap_{n \geq 0} \mathcal{F}_{n,\infty}^X \vee \mathcal{F}_{0,\infty}^Y}(\omega) \\ &\quad - \int_A f_2(\omega) \log(f_2(\omega)) dP^\nu|_{\mathcal{F}_{0,\infty}^Y}(\omega) \\ &= \int_A f_1(\omega) \log(f_1(\omega)) dP^\nu|_{\mathcal{F}_{0,\infty}^Y}(\omega) \\ &\quad - \int_A f_2(\omega) \log(f_2(\omega)) dP^\nu|_{\mathcal{F}_{0,\infty}^Y}(\omega) = 0 \end{aligned}$$

Therefore, if the relative entropy of the filter is ever finite, then total variation merging in expectation is equivalent to merging in relative entropy. \square

V. CONCLUSION

We presented a notion of stochastic observability for non-linear systems. This notion is explicit, is relatively easily computed due to its functional approximation formulation, and is shown via examples to be applicable to a large class of systems. The implications of this definition for filter stability were presented in detail. Further relations under various stability criteria and implications were studied.

VI. APPENDIX. SUPPORTING RESULTS FOR SECTION IV-E

We present a number of supporting results. The approach for these build on similar arguments in [14] and [54]. The proofs here are kept brief due to space constraints or omitted.

Lemma VI.1. Assume $\mu \ll \nu$. For any sigma field $\mathcal{G} \subseteq \mathcal{F}_{0,\infty}^X \vee \mathcal{F}_{0,\infty}^Y$ we have:

$$\frac{dP^\mu|_{\mathcal{G}}}{dP^\nu|_{\mathcal{G}}} = E^\nu \left[\frac{d\mu}{d\nu}(X_0) \Big| \mathcal{G} \right] \quad P^\mu \text{ a.s.}$$

Lemma VI.2. Assume $\mu \ll \nu$. For any two sigma fields $\mathcal{G}_1, \mathcal{G}_2 \subset \mathcal{F}_{0,\infty}^X \vee \mathcal{F}_{0,\infty}^Y$, let $P^\mu|_{\mathcal{G}_1} | \mathcal{G}_2$ represent the probability measure P^μ restricted to \mathcal{G}_1 , conditioned on field \mathcal{G}_2 . We then have

$$\frac{dP^\mu|_{\mathcal{G}_1} | \mathcal{G}_2}{dP^\nu|_{\mathcal{G}_1} | \mathcal{G}_2} = \frac{E^\nu \left[\frac{d\mu}{d\nu}(X_0) | \mathcal{G}_1 \vee \mathcal{G}_2 \right]}{E^\nu \left[\frac{d\mu}{d\nu}(X_0) | \mathcal{G}_2 \right]} \quad P^\mu \text{ a.s.}$$

Lemma VI.3. Assume $\mu \ll \nu$, for any two sigma fields $\mathcal{G}_1, \mathcal{G}_2 \subset \mathcal{F}_{0,\infty}^{\mathcal{X}} \vee \mathcal{F}_{0,\infty}^{\mathcal{Y}}$ we have P^μ a.s.

$$\begin{aligned} & \|P^\mu|_{\mathcal{G}_1} \mathcal{G}_2 - P^\nu|_{\mathcal{G}_1} \mathcal{G}_2\|_{TV} \\ &= \frac{E^\nu \left[\left| E^\nu \left[\frac{d\mu}{d\nu}(X_0) | \mathcal{G}_1 \vee \mathcal{G}_2 \right] - E^\nu \left[\frac{d\mu}{d\nu}(X_0) | \mathcal{G}_2 \right] \right| \middle| \mathcal{G}_2 \right]}{E^\nu \left[\frac{d\mu}{d\nu}(X_0) | \mathcal{G}_2 \right]} \end{aligned}$$

For the specific case of the non-linear filter, that is $\mathcal{G}_1 = \mathcal{F}_n^{\mathcal{X}}$ and $\mathcal{G}_2 = \mathcal{F}_{0,n}^{\mathcal{Y}}$, the results presented above imply the following known results in the literature.

Lemma VI.4. [54, Lemma 5.6] Assume $\mu \ll \nu$. Then we have that $\pi_n^\mu \ll \pi_n^\nu$ a.s. and we have

$$\frac{d\pi_n^\mu(x)}{d\pi_n^\nu(x)} = \frac{E^\nu \left[\frac{d\mu}{d\nu}(X_0) | Y_{[0,n]}, X_n = x \right]}{E^\nu \left[\frac{d\mu}{d\nu}(X_0) | Y_{[0,n]} \right]} \quad P^\mu \text{ a.s.} \quad (15)$$

Lemma VI.5. [14, Equation 1.10] The filter merges in total variation in expectation if and only if P^ν a.s.

$$E^\nu \left[\frac{d\mu}{d\nu}(X_0) \middle| \bigcap_{n \geq 0} \mathcal{F}_{0,\infty}^{\mathcal{Y}} \vee \mathcal{F}_{n,\infty}^{\mathcal{X}} \right] = E^\nu \left[\frac{d\mu}{d\nu}(X_0) \middle| \mathcal{F}_{0,\infty}^{\mathcal{Y}} \right] \quad (16)$$

Since our results apply to any general sigma field, not just the fields used in the analysis of the filter, we can study the predictor process to establish Lemmas VI.6, VI.7, and VI.8, in the following.

Lemma VI.6. Assume $\mu \ll \nu$. Then we have that $\pi_{n-}^\mu \ll \pi_{n-}^\nu$ P^μ a.s. and we have

$$\frac{d\pi_{n-}^\mu(x)}{d\pi_{n-}^\nu(x)} = \frac{E^\nu \left[\frac{d\mu}{d\nu}(X_0) | Y_{[0,n-1]}, X_n = x \right]}{E^\nu \left[\frac{d\mu}{d\nu}(X_0) | Y_{[0,n-1]} \right]} \quad P^\mu \text{ a.s.} \quad (17)$$

Proof. These results become clear from Lemma VI.2 when we state the predictor as P^μ restricted to $\mathcal{F}_n^{\mathcal{X}}$ conditioned on $\mathcal{F}_{0,n-1}^{\mathcal{Y}}$ \square

Lemma VI.7. Assume $\mu \ll \gamma$ for some measure γ . We can express

$$\begin{aligned} & \|\pi_{n-}^\mu - \pi_{n-}^\gamma\|_{TV} = \\ & \frac{E^\gamma \left[\left| E^\gamma \left[\frac{d\mu}{d\gamma}(X_0) | Y_{[0,\infty)}, X_{[n,\infty)} \right] - E^\gamma \left[\frac{d\mu}{d\gamma}(X_0) | Y_{[0,n-1]} \right] \right| \middle| Y_{[0,n-1]} \right]}{E^\gamma \left[\frac{d\mu}{d\gamma}(X_0) \middle| Y_{[0,n-1]} \right]} \end{aligned} \quad (18)$$

Proof. By Lemma VI.3 we can write

$$\begin{aligned} & \|\pi_{n-}^\mu - \pi_{n-}^\gamma\|_{TV} = \\ & \frac{E^\gamma \left[\left| E^\gamma \left[\frac{d\mu}{d\gamma}(X_0) | Y_{[0,n-1]}, X_n \right] - E^\gamma \left[\frac{d\mu}{d\gamma}(X_0) | Y_{[0,n-1]} \right] \right| \middle| Y_{[0,n-1]} \right]}{E^\gamma \left[\frac{d\mu}{d\gamma}(X_0) \middle| Y_{[0,n-1]} \right]} \end{aligned}$$

Since Y_n is a function of X_n and the random noise Z_n which is independent of X_n and past $Y_{[0,n-1]}$ measurements, we have that $\sigma(Y_{[0,n-1]}, X_n) = \sigma(Y_{[0,n]}, X_n)$. Further, by the Markov property we have that we have that

$(X_{[0,n-1]}, Y_{[0,n-1]})$ are independent of $(X_{[n+1,\infty)}, Y_{[n+1,\infty)})$ conditioned on (X_n, Y_n) therefore we can state

$$E^\gamma \left[\frac{d\mu}{d\nu}(X_0) | Y_{[0,n-1]}, X_n \right] = E^\gamma \left[\frac{d\mu}{d\nu}(X_0) | Y_{[0,\infty)}, X_{[n,\infty)} \right] \quad \square$$

Lemma VI.8. The predictor merges in total variation in expectation if and only if

$$E^\nu \left[\frac{d\mu}{d\nu}(X_0) \middle| \bigcap_{n \geq 1} \mathcal{F}_{0,\infty}^{\mathcal{Y}} \vee \mathcal{F}_{n,\infty}^{\mathcal{X}} \right] = E^\nu \left[\frac{d\mu}{d\nu}(X_0) \middle| \mathcal{F}_{0,\infty}^{\mathcal{Y}} \right] \quad P^\nu \text{ a.s.} \quad (19)$$

Proof. Building on the proof of Lemma VI.7, we have

$$\begin{aligned} E^\mu \left[\|\pi_{n-}^\mu - \pi_{n-}^\nu\|_{TV} \right] &= E^\nu \left[\frac{dP^\mu|_{\mathcal{F}_{0,n-1}^{\mathcal{Y}}}}{dP^\nu|_{\mathcal{F}_{0,n-1}^{\mathcal{Y}}}} \|\pi_{n-}^\mu - \pi_{n-}^\nu\|_{TV} \right] \\ &= E^\nu \left[E^\nu \left[\frac{d\mu}{d\nu}(X_0) \middle| Y_{[0,n-1]} \right] \|\pi_{n-}^\mu - \pi_{n-}^\nu\|_{TV} \right] \\ &= E^\nu \left[E^\nu \left[E^\nu \left[\frac{d\mu}{d\nu}(X_0) | Y_{[0,\infty)}, X_{[n,\infty)} \right] \right. \right. \\ &\quad \left. \left. - E^\nu \left[\frac{d\mu}{d\nu}(X_0) | Y_{[0,n-1]} \right] \middle| Y_{[0,n-1]} \right] \right] \\ &= E^\nu \left[\left| E^\nu \left[\frac{d\mu}{d\nu}(X_0) | Y_{[0,\infty)}, X_{[n,\infty)} \right] - E^\nu \left[\frac{d\mu}{d\nu}(X_0) | Y_{[0,n-1]} \right] \right| \right] \end{aligned}$$

We then see that $A_n = E^\nu \left[\frac{d\mu}{d\nu}(X_0) | Y_{[0,n-1]} \right]$ is a non-negative uniformly integrable martingale adapted to the increasing filtration $\mathcal{F}_{0,n-1}^{\mathcal{Y}}$. Hence the limit as $n \rightarrow \infty$ in $L^1(P^\nu)$ is $E^\nu \left[\frac{d\mu}{d\nu}(X_0) | \mathcal{F}_{0,\infty}^{\mathcal{Y}} \right]$. Similarly, we can view $B_n = E^\nu \left[\frac{d\mu}{d\nu}(X_0) | Y_{[0,\infty)}, X_{[n,\infty)} \right]$ as backwards non-negative uniformly integrable martingale with respect to the decreasing sequence of filtrations $\mathcal{F}_{0,\infty}^{\mathcal{Y}} \vee \mathcal{F}_{n,\infty}^{\mathcal{X}}$. Then by the backwards martingale convergence theorem, the limit as $n \rightarrow \infty$ in $L^1(P^\nu)$ is $E^\nu \left[\frac{d\mu}{d\nu}(X_0) | \bigcap_{n=1}^\infty \mathcal{F}_{1,\infty}^{\mathcal{Y}} \vee \mathcal{F}_{n,\infty}^{\mathcal{X}} \right]$. It is then clear the total variation in expectation is zero if and only if equation (19) holds. \square

REFERENCES

- [1] B. D. O. Anderson and J. B. Moore. Detectability and stabilizability of time-varying discrete-time linear systems. *SIAM Journal on Control and Optimization*, 19(1):20–32, 1981.
- [2] Y. Bar-Shalom and E. Tse. Dual effect, certainty equivalence, and separation in stochastic control. *IEEE Transactions on Automatic Control*, 19(5):494–500, October 1974.
- [3] A. R. Barron. Information theory and martingales. In *IEEE International Symposium on Information Theory, recent results session*, Budapest, Hungary, 1991.
- [4] A. R. Barron. Limits of information, Markov chains, and projections. Sorrento, Italy, 2000. Proceedings of the IEEE Int. Symp. on Inform. Theory p. 25.
- [5] P. Billingsley. *Probability and Measure*. Wiley, New York, 2nd edition, 1986.
- [6] P. Billingsley. *Convergence of probability measures*. New York: Wiley, 2nd edition, 1999.
- [7] D. Blackwell and L. Dubins. Merging of opinions with increasing information. *The Annals of Mathematical Statistics*, 33(3):882–886, 1962.
- [8] V. I. Bogachev. *Measure Theory*. Springer-Verlag, Berlin, 2007.
- [9] V. S. Borkar. White-noise representations in stochastic realization theory. *SIAM J. on Control and Optimization*, 31:1093–1102, 1993.
- [10] P. E. Caines. *Linear Stochastic Systems*. John Wiley & Sons, New York, NY, 1988.

- [11] C. T. Chen. *Linear Systems Theory and Design*. Oxford University Press, Oxford, 1999.
- [12] P. Chigansky and R. Van Handel. A complete solution to Blackwell's unique ergodicity problem for hidden Markov chains. *The Annals of Applied Probability*, 20(6):2318–2345, 2010.
- [13] P. Chigansky and R. Liptser. On a role of predictor in the filtering stability. *Electron. Comm. Probab.*, 11:129–140, 2006.
- [14] P. Chigansky, R. Liptser, and R. van Handel. Intrinsic methods in filter stability. *Handbook of Nonlinear Filtering*, 2009.
- [15] J.M.C. Clark, D. L. Ocone, and C. Coumarbatch. Relative entropy and error bounds for filtering of markov processes. *Mathematics of Control, Signals and Systems*, 12(4):346–360, 1999.
- [16] I. Csiszár. Information-type measures of difference of probability distributions and indirect observation. *studia scientiarum Mathematicarum Hungarica*, 2:229–318, 1967.
- [17] A. D'Aristotile, P. Diaconis, and D. Freedman. On merging of probabilities. *Sankhyā: The Indian Journal of Statistics, Series A*, pages 363–380, 1988.
- [18] R. M. Dudley. *Real Analysis and Probability*. Cambridge University Press, Cambridge, 2nd edition, 2002.
- [19] J. Dugundji. An extension of tietze's theorem. *Pacific Journal of Mathematics*, 1(3):353–367, 1951.
- [20] S. N. Ethier and T. G. Kurtz. *Markov processes: characterization and convergence*, volume 282. John Wiley & Sons, 2009.
- [21] E.A. Feinberg, P.O. Kasyanov, and M.Z. Zgurovsky. Partially observable total-cost Markov decision process with weakly continuous transition probabilities. *Mathematics of Operations Research*, 41(2):656–681, 2016.
- [22] J. Fritz. An information-theoretical proof of limit theorems for reversible Markov processes. In *In Trans. Sixth Prague Conf. on Inform. Theory, Statist. Decision Functions, Random Processes*, 1973.
- [23] I. I. Gihman and A. V. Skorohod. *Controlled stochastic processes*. Springer Science & Business Media, 2012.
- [24] N. Golowich, A. Moitra, and D. Rohatgi. Planning in observable pomdps in quasipolynomial time. *arXiv preprint arXiv:2201.04735*, 2022.
- [25] R. M. Gray. *Entropy and information theory*. Springer Science & Business Media, 2011.
- [26] P. Harremoës and K. K. Holst. Convergence of Markov chains in information divergence. *J. Theoretical Probability*, 22:186–202, 2009.
- [27] P. Harremoës, O. Johnson, and I. Kontoyiannis. Thinning, entropy, and the law of thin numbers. *IEEE Transactions on Information Theory*, 56(9):4228–4244, 2010.
- [28] R. Hermann and A. Krener. Nonlinear controllability and observability. *IEEE Transactions on Automatic Control*, 22(5):728–740, 1977.
- [29] C. C. White III and W. T. Scherer. Finite-memory suboptimal design for partially observed markov decision processes. *Operations Research*, 42(3):439–455, 1994.
- [30] R. E. Kalman. A new approach to linear filtering and prediction problems. *Transaction of the ASME J. of Basic Engineering*, pages 35–45, March 1960.
- [31] A.D Kara, N. Saldi, and S. Yüksel. Weak Feller property of non-linear filters. *Systems & Control Letters*, 134:104–512, 2019.
- [32] A.D Kara and S. Yüksel. Robustness to incorrect priors in partially observed stochastic control. *SIAM Journal on Control and Optimization*, 57(3):1929–1964, 2019.
- [33] A.D Kara and S. Yüksel. Convergence of finite memory Q-learning for POMDPs and near optimality of learned policies under filter stability. *Mathematics of Operations Research (also arXiv:2103.12158)*, 2022.
- [34] A.D Kara and S. Yüksel. Near optimality of finite memory feedback policies in partially observed markov decision processes. *Journal of Machine Learning Research*, 23(11):1–46, 2022.
- [35] J.-W. Kim and P. G. Mehta. Duality for nonlinear filtering i: Observability. *arXiv preprint arXiv:2208.06586*, 2022.
- [36] J.-W. Kim, P. G. Mehta, and S. P. Meyn. What is the Lagrangian for nonlinear filtering? In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 1607–1614. IEEE, 2019.
- [37] H. Kunita. Asymptotic behavior of the nonlinear filtering errors of markov processes. *Journal of Multivariate Analysis*, 1(4):365–393, 1971.
- [38] H. J. Kushner. A partial history of the early development of continuous-time nonlinear stochastic systems theory. *Automatica*, 50(2):303–334, 2014.
- [39] H.J. Kushner. *Introduction to Stochastic Control Theory*. Holt, Rinehart and Winston, New York, 1972.
- [40] T. Linder and S. Yüksel. On optimal zero-delay coding of vector markov sources. *IEEE Trans. Information Theory*, 60(10):5975–5991, 2014.
- [41] A. R. Liu. *Stochastic observability, reconstructibility, controllability, and reachability*. PhD thesis, UC San Diego, 2011.
- [42] A.R. Liu and R.R. Bitmead. Observability and reconstructibility of hidden markov models: Implications for control and network congestion control. In *49th IEEE Conference on Decision and Control (CDC)*, pages 918–923. IEEE, 2010.
- [43] C. McDonald and S. Yüksel. Robustness to incorrect priors and controlled filter stability in partially observed stochastic control. *SIAM Journal on Control and Optimization*, 60(2):842–870, 2022.
- [44] S. P. Meyn and R. Tweedie. *Markov Chains and Stochastic Stability*. Springer-Verlag, London, 1993.
- [45] H. Nijmeijer. Observability of autonomous discrete time non-linear systems: a geometric approach. *International journal of control*, 36(5):867–874, 1982.
- [46] D. L. Ocone. Asymptotic stability of beneš filters. *Stochastic analysis and applications*, 17(6):1053–1074, 1999.
- [47] A.S. Reddy and A. Apte. Stability of non-linear filter for deterministic dynamics. *arXiv preprint arXiv:1910.14348*, 2019.
- [48] W. Rudin. *Real and Complex Analysis*. Tata McGraw-Hill Education, 2006.
- [49] E. D. Sontag. A concept of local observability. *Systems & Control Letters*, 5(1):41–47, 1984.
- [50] R.B Sowers and A.M. Makowski. Discrete-time filtering for linear systems with non-Gaussian initial conditions: asymptotic behavior of the difference between the MMSE and LMSE estimates. *IEEE transactions on automatic control*, 37(1):114–120, 1992.
- [51] V.A. Ugrinovskii. Observability of linear stochastic uncertain systems. *IEEE Transactions on Automatic Control*, 48(12):2264–2269, 2003.
- [52] R. van Handel. Discrete time nonlinear filters with informative observations are stable. *Electronic Communications in Probability*, 13:562–575, 2008.
- [53] R. van Handel. Observability and nonlinear filtering. *Probability theory and related fields*, 145(1-2):35–74, 2009.
- [54] R. van Handel. The stability of conditional markov processes and markov chains in random environments. *The Annals of Probability*, 37(5):1876–1925, 2009.
- [55] R. van Handel. Uniform observability of hidden Markov models and filter stability for unstable signals. *The Annals of Applied Probability*, 19(3):1172–1199, 2009.
- [56] R. van Handel. Nonlinear filtering and systems theory. In *Proceedings of the 19th International Symposium on Mathematical Theory of Networks and Systems (MTNS semi-plenary paper)*, 2010.