

Max-Planck-Institut
für Mathematik
in den Naturwissenschaften
Leipzig

On the Generative Nature of Prediction

(revised version: September 2008)

by

Wolfgang Löhr, and Nihat Ay

Preprint no.: 8

2008



On the Generative Nature of Prediction

Wolfgang Löhr¹ & Nihat Ay^{1,2}

August 14, 2008

Abstract

Given an observed stochastic process, computational mechanics provides an explicit and efficient method of constructing a minimal hidden Markov model within the class of maximally predictive models. Here, the corresponding so-called ε -machine encodes the mechanisms of prediction. We propose an alternative notion of predictive models in terms of a hidden Markov model capable of generating the underlying stochastic process. A comparison of these two notions of prediction reveals that our approach is less restrictive and thereby allows for predictive models that are more concise than the ε -machine.

Keywords: hidden Markov models, computational mechanics, ε -machines, prediction

Contents

1	Introduction	2
1.1	Models of prediction and their constraints	2
1.2	The main idea of the paper	3
1.3	Generative models: Stochastic automata, HMMs, and OOMs	7
2	Predictive models of stochastic processes	8
2.1	Generating a process	8
2.2	Our prediction setting	9
3	What does “prediction” really mean?	13
3.1	Predictive versus prescient memories	13
3.2	Implications on minimality of predictive models	16
4	Conclusions	19
A	Appendix: Prediction with infinite histories	19
A.1	Generating a process	20
A.2	Our prediction setting	20
A.3	Comparison to computational mechanics	21

¹Max Planck Institute for Mathematics in the Sciences, Inselstraße 22, D-04103 Leipzig, Germany

²Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, New Mexico 87501, USA

1 Introduction

1.1 Models of prediction and their constraints

This paper is mainly about computational mechanics, a theory introduced and further developed by Crutchfield and coworkers [Crutchfield and Young, 1989, Shalizi and Crutchfield, 2001, Ay and Crutchfield, 2005, Still and Crutchfield, 2007]. It deals with the following problems: Having an observed stochastic process, what is the best model in the sense of minimal size and maximal predictive power? If such a model exists, is there a way of explicitly constructing it? How much memory is needed to generate the process? In the context of these problems, the so-called ε -machine and its construction have been proposed as the optimal solution.

The models considered in computational mechanics are stochastic output automata, also called edge-emitting hidden Markov models (HMMs). We give a short discussion of the different types of HMMs in Section 1.3. The ε -machine is directly derived from the so-called causal states which are the minimal sufficient statistic on the past for predicting the future of the stochastic process. More precisely, they form the coarsest partition of the past that retains all information about the future. Furthermore, they can be obtained in a constructive and efficient way ([Shalizi et al., 2002]).

There may exist, however, HMMs with fewer internal states than the corresponding ε -machine. There are even examples of HMMs with few internal states that lead to an ε -machine with infinite size ([Crutchfield, 1994], see also Example 3.6). Consequently, ε -machines can only be of minimal size within reduced model classes that satisfy appropriately chosen constraints. Unfortunately, in the literature these constraints are mainly implicitly stated, which leads to a misperception and the above-mentioned apparent inconsistency. In particular, the ε -machine is not the minimal model that *generates* a given process, a property occasionally referred to in the literature. Instead, the ε -machine minimality depends on constraints that are related to particular notions of memory and prediction based on the fundamental concept of sufficient statistic. This paper reveals and relaxes these constraints in such a way that, given an HMM, a minimal predictive model of the corresponding output process can not be larger than that HMM. Recently, one step towards such an extension, namely introducing stochastic memories, has been made in [Still and Crutchfield, 2007]. While the focus of that paper lies on the trade-off between predictive power and model size (allowing some prediction error) based on the bottleneck method [Tishby et al., 1999], we show in the present contribution that this extension also is necessary for having concise models with maximal predictive power (in our weakened sense). On the other hand, it allows for smaller models only in combination with our notion of prediction.

In Section 1.2, we provide an intuitive sketch of our main idea in a simplified setting before going into the technical details of precise definitions and results within the main part of the paper. Based on our understanding of prediction as a generative operation we then briefly discuss several generative models for stochastic processes in Section 1.3. The main part, Sections 2 and 3, contains definitions and results for the case of prediction with finite histories and does not require any measure theoretic background. To complete the discussion, in the appendix we apply measure theoretic tools for showing that all results, up to minor modifications, remain true in the cases of infinite history and general state spaces.

1.2 The main idea of the paper

In order to illustrate the main idea of the paper, as an initial step we start with a simplified situation of a one-step prediction scenario. More precisely, we consider a pair X_p and X_f of discrete random variables, which we interpret as past and future observations. Not all information of X_p is necessary for predicting X_f , so that one tries to compress the relevant information in a (discrete) memory variable M via a memory map mem as shown in Figure 1.

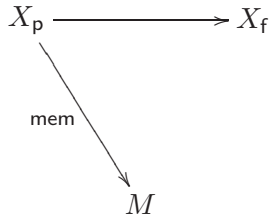


Figure 1: Memory map mem that compresses the information contained in X_p about X_f

Within computational mechanics, mainly deterministic memory maps mem have been considered, where M is assumed to be a *sufficient statistic*¹ on the past X_p for the future X_f . Only recently an extension to stochastic maps has been considered ([Still and Crutchfield, 2007]). We adopt this extension and do not require mem to be a deterministic function but allow for a stochastic assignment, i.e. the memory map is assumed to be a Markov kernel (transition probability).

In general, the map mem reduces the information about the future, which is expressed by the following inequality:

$$I(M : X_f) \leq I(X_p : X_f), \quad (1)$$

where I denotes the mutual information between two variables. The mutual information $I(X_p : X_f)$ between past and future corresponds to a complexity measure defined within the context of stochastic processes and known as *effective measure complexity*, *excess entropy*, and *predictive information* [Grassberger, 1986, Shalizi and Crutchfield, 2001, Bialek et al., 2001]. We use the term predictive information for $I(X_p : X_f)$ also within the simplified setting of this section.

In computational mechanics one assumes that M is *prescient* ([Shalizi and Crutchfield, 2001]) in the sense that the equality holds in (1):

$$I(M : X_f) = I(X_p : X_f). \quad (2)$$

Note that for a deterministic memory map mem , this property reduces to the notion of sufficient statistic. Furthermore, (2) directly implies that the predictive information is a lower bound of the entropy $H(M)$ which corresponds to a complexity measure known as *statistical complexity* within computational mechanics ([Shalizi and Crutchfield, 2001]):

$$I(X_p : X_f) \leq H(M) \leq \ln(|M|), \quad (3)$$

¹See, e.g., [Kulhavý, 1996] for an introduction to statistics in an abstract context.

where $|M|$ denotes the cardinality of the set M of memory states.

The requirement (2) is equivalent to

$$I(X_f : X_p | M) = 0, \quad (4)$$

which means that M captures all the information of the past that is necessary for predicting the future, or, stated differently, the future is independent of the past given the memory:

$$X_f \perp\!\!\!\perp X_p | M. \quad (5)$$

This implies that the joint distribution can be represented by the graph of Figure 2.

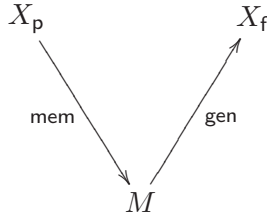


Figure 2: Factorization of the channel $X_p \rightarrow X_f$ resulting from a prescient memory

More precisely, there is a Markov kernel gen , illustrated by the arrow $M \rightarrow X_f$, that satisfies

$$\mathbb{P}(X_f = x_f | X_p = x_p) = \sum_m \text{mem}(x_p; m) \text{gen}(m; x_f). \quad (6)$$

Interpreting gen in a generative way, this means that one can replace the mechanisms that underly the channel $X_p \rightarrow X_f$ by the two steps of

1. mapping the past to a memory state using mem and then, based on that state,
2. generating the future with gen .

It is important to point out that the replacement of the channel $X_p \rightarrow X_f$ by the composition of mem and gen does not mean that both mechanisms have the same outcome. Due to the intrinsic stochasticity of the kernels, one can not expect that the predicted outcome coincides with the real outcome. Only the distribution is the same. In order to explicitly model this outcome difference, we use the symbol \tilde{X}_f to denote the prediction variable, given by the mechanisms mem and gen , and combine the graphs of the Figures 1 and 2 into the graph of Figure 3. In this picture, the condition (6) translates to the following condition:

$$\mathbb{P}(X_f = x_f | X_p = x_p) = \mathbb{P}(\tilde{X}_f = x_f | X_p = x_p). \quad (7)$$

Here, the existence of a kernel gen that satisfies (7) is a consequence of the assumption that the memory is prescient. On the other hand, one could try to find a kernel mem , not necessarily prescient, for which there exists a kernel gen such that (7) is satisfied. The requirement (7) simply means that the memory mem contains sufficient information for generating a future trajectory with gen that is indistinguishable from the real future trajectory based on the observed past. We call such a memory *predictive*. This is the ansatz of the present paper,

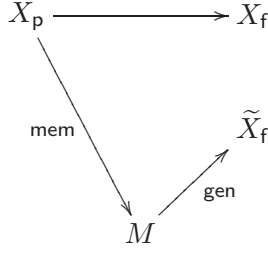


Figure 3: Memory map `mem` together with a generator `gen` which generates \tilde{X}_f as a version of X_f that is indistinguishable from X_f based on X_p

which relaxes the notion of a prescient memory. At the end of this section, we give an example of a predictive memory that is not prescient.

With inequality (1) and with the assumption (7) one has the following relations to the predictive information, that is the mutual information between past and future:

$$I(M : X_f) \leq I(X_p : X_f) = I(X_p : \tilde{X}_f) \leq I(M : \tilde{X}_f). \quad (8)$$

More precisely, only the equality in (8) requires predictability of the memory. According to the second inequality, a memory can contain more information about the predicted future \tilde{X}_f than the past X_p , whereas, according to the first inequality, no memory can contain more information about the actual future than the past. In view of the usual reference to the data processing inequality within the computational mechanics literature, the second inequality appears unfamiliar and emphasizes a new aspect of prediction. On the other hand, it is simply a direct consequence of the explicit distinction between the actual and predicted future, that is X_f and \tilde{X}_f . Furthermore, as a direct implication of (8) the inequality (3), which holds for prescient memories, remains valid also for predictive memories. Stated explicitly, the entropy of a predictive memory variable M is lower bounded by the predictive information of the observed process. We will see that, as consequence of our extension of the memory class to predictive memories, it is possible to come closer to this lower bound and thereby better exploit the predictive structure of the underlying process.

Example 1.1 shows the difference between prescient and predictive memories within the simplified setting of this introduction. As we already mentioned, all definitions of this section will be extended to the general setting of stochastic processes. The key idea of our main Example 3.6 within that general setting is illustrated by the following construction:

Example 1.1. Let X_p and X_f assume values in $\{0, 1, 2\}$,

$$\mathbb{P}(X_p = x_p) = \frac{1}{3} \quad \text{and} \quad \mathbb{P}(X_f = x_f | X_p = x_p) = \begin{cases} \frac{2}{3}, & \text{if } x_f = x_p \neq 2 \\ \frac{1}{3}, & \text{if } x_p = 2 \text{ or } x_f = 2. \\ 0, & \text{otherwise} \end{cases}$$

As predictive information we get

$$I(X_p : X_f) = \frac{4}{9} \ln(2),$$

which is a lower bound for the entropy of any prescient or predictive memory. Since all three possibilities for X_p yield different expectations on the future, every prescient memory variable still has to distinguish between these three states (the mathematical proof of this statement is simple and contained in the proof of Proposition 3.5). Intuitively, if we identify states x_p with the same conditional distribution on the x_f we get three so-called causal states as equivalence classes. Obviously, every prescient memory has at least three memory states with an entropy $H(M) \geq \ln(3)$. On the other hand, it is possible to define a predictive memory with only two memory states. To this end, consider a memory with state set $\{0, 1\}$ that copies x_p if $x_p \in \{0, 1\}$ and chooses a random state otherwise:

$$\mathbb{P}(M = m \mid X_p = x_p) = \text{mem}(x_p; m) := \begin{cases} 1, & \text{if } m = x_p \neq 2 \\ \frac{1}{2}, & \text{if } x_p = 2 \\ 0, & \text{otherwise} \end{cases}.$$

We see that mem is predictive by defining $\text{gen}(m) = \mathbb{P}(X_f \mid X_p = m)$, $m \in \{0, 1\}$ and observing that also for $x_p = 2$ this yields the correct distribution of the prediction:

$$\begin{aligned} \mathbb{P}(\tilde{X}_f \mid X_p = 2) &= \frac{1}{2}\mathbb{P}(\tilde{X}_f \mid M = 0) + \frac{1}{2}\mathbb{P}(\tilde{X}_f \mid M = 1) \\ &= \frac{1}{2}\mathbb{P}(X_f \mid X_p = 0) + \frac{1}{2}\mathbb{P}(X_f \mid X_p = 1) = \mathbb{P}(X_f \mid X_p = 2). \end{aligned}$$

Furthermore, $H(M) = \ln(2)$. ◇

Within the general context of stochastic processes, our main results are the following: Every prescient memory map mem , which is assumed to satisfy (2), is also predictive but not vice versa (Proposition 3.4 and Example 3.6). On the other hand, restricted to the situation where the map mem is deterministic, i.e. M is a statistic on X_p , predictive and prescient turn out to be equivalent (Proposition 3.4). The ε -machine turns out to be minimal in the class of *prescient* memories, even if we allow stochastic memory maps (Proposition 3.5). Every HMM, however, can be considered as a predictive model with the same number of memory states (Proposition 2.5), which provides a way of constructing machines that are more concise than the ε -machine. Example 3.6 is such a (minimal, edge-emitting) HMM with two hidden states whereas the corresponding ε -machine requires infinitely many hidden states.

Although extending the model class from prescient to predictive allows for substantially smaller models, while preserving a notion of predictive power which we consider quite natural, we have to point out two drawbacks. Firstly, constructing a minimal HMM is intrinsically difficult, whereas efficient algorithms are available for the construction of the ε -machine. Secondly and conceptually more important, the property of sufficiency is lost. This means that the memory state is not a complete substitute for the past. Within computational mechanics, the total information that is required for encoding the future distribution is completely contained in the memory state whereas in our model the memory only contains that part of it that is sufficient for generating the future distribution based on the past. When we want to use our model for sampling the future distribution given an x_p , we first choose a memory state m according to $\text{mem}(x_p)$ and then apply the process $\text{gen}(m)$ which generates an x_f . If we repeat this sampling procedure we obtain the correct future distribution. On the other hand, if we “forget” the history state x_p and, instead of sampling new m ’s according to $\text{mem}(x_p)$, apply $\text{gen}(m)$ using the same m , the resulting distribution of x_f will be different from $\mathbb{P}(X_f \mid X_p = x_p)$. Thus, we have to memorize the *distribution* (the information state) of

the memory states m given x_p . One can show that the number of these information states is lower bounded by the number of causal states.

1.3 Generative models: Stochastic automata, HMMs, and OOMs

In computational mechanics the ε -machine is defined in terms of a *stochastic (output) automaton*.² This name is directly linked to the intuition of a “machine”: It has internal states M and is initialized by one of these states according to some initial probability distribution. At each time step t the internal state M_t is updated, that is M_{t+1} is generated and, at the same time, an output symbol X_{t+1} from a finite alphabet D is emitted. This is modeled by a joint transition probability from the internal states to output and internal states. Stochastic automata are also widely known as *edge-emitting hidden Markov models*. Here “edge-emitting” means that the output symbol may depend on both the old and the new internal state (output symbol and new internal state are determined by a joint kernel). In transition graph representations of edge-emitting HMMs, the output symbols appear as edge labels.

Probably more common than edge-emitting HMMs are the more restrictive *state-emitting HMMs*. These HMMs have to satisfy the additional condition that the output symbol X_{t+1} depends on either M_t or M_{t+1} but not on both (as in the case of edge-emitting HMMs). More precisely, it is assumed that the transition probability factorizes as follows:

$$\mathbb{P}(M_{t+1}, X_{t+1} | M_t) = \mathbb{P}(M_{t+1} | M_t) \mathbb{P}(X_{t+1} | M_s), \quad s \in \{t, t+1\}.$$

Even more restrictive are *functions of Markov chains*. Here, the output symbol is a deterministic function of the internal state M_{t+1} . Sometimes also a notion less restrictive than edge-emitting HMM is considered, namely that of a *partially observed Markov process*:³ A (time homogeneous) Markov process $(M_t, X_t)_{t \in \mathbb{Z}}$ on a product space $M \times D$, where the component D is considered to be observable, whereas the other component M consists of hidden states. Note that here *both* marginal processes $(M_t)_{t \in \mathbb{Z}}$ and $(X_t)_{t \in \mathbb{Z}}$ need not be Markovian.

These four notions, partially observed Markov process, edge-emitting HMM (stochastic automaton), state-emitting HMM, and function of a Markov chain (ordered from more general to more restrictive) are essentially equivalent in the following sense: To every partially observed Markov process, one can naturally associate a function of a Markov chain such that the cardinality of the internal states increases only by the constant factor of the cardinality of the output alphabet. One simply takes as new set M' of internal states the product $M \times D$ of internal and observable states and the projection $m' = (m, x) \mapsto x$ onto the observable component x as function determining the output symbol. In the following, we only consider edge-emitting HMMs and HMM always means edge-emitting.

There are also more algebraic models of stochastic processes, which dismiss the conception of the internal dynamics being described by a stochastic process. Instead, some vector space replaces the internal states and the “dynamics” is described by linear maps (instead of Markov kernels). These models were introduced and termed *stochastic S -modules* by Heller in the very concise and well-written paper [Heller, 1965]. Later, in [Jaeger, 2000], Jaeger made the construction more explicit and transparent for readers not familiar with module theory. He

²See, e.g., [Bukharaev, 1995] for an introduction to the theory of stochastic automata.

³The term “partially observed Markov process” sometimes refers to edge- or state-emitting HMMs.

introduced the name *observable operator model (OOM)*, provided ways of interpreting them and extended the theory by learning algorithms. Ergodic theory for OOMs was developed in [Faigle and Schönhuth, 2007, Schönhuth and Jaeger, 2007]. In [Littman et al., 2001], the same model class was also obtained starting from a somewhat different intuition (internal states are constructed as predictions for certain tests) as linear non-controlled *predictive state representations (PSRs)*.⁴ The equivalence of linear non-controlled PSRs and OOMs is shown in [Singh et al., 2004].

OOMs are a generalization of HMMs in the sense that for any HMM there is a naturally associated, equivalent OOM with the internal states of the HMM providing a basis for its internal vector space.⁵ Here the OOM state vectors roughly correspond to the so called *information states* (probability distributions over internal states) and the number of internal states of the HMM is equal to the dimension of the associated OOM. However, not every OOM is induced by an HMM (as the transitions need not be positive), and the dimension of the minimal OOM may be substantially smaller than the minimal number of HMM states. In fact, finite dimensional OOMs exist for some processes which do not allow for HMMs with finitely many internal states.

While OOMs can be constructed such that the internal state vector has some probabilistic interpretation (giving the probability of certain future events), it may be impossible to interpret any basis as “internal (pure) states”: To simulate a process with an HMM, one can, at each time step, determine the next internal state randomly. Thus one has to store the internal state only (as opposed to the information state). This is not possible for OOMs in general, due to the occurrence of “negative probabilities”.

The question whether there exists a (finite) HMM, as well as finding the minimal one, is intrinsically difficult. It depends on an intricate geometrical condition specified in [Heller, 1965]. Thus, although we want to keep the interpretation of a probabilistic machine and work with HMMs as generative models, OOM theory is interesting in that it provides a constructive algorithm for finding the minimal OOM. In some cases, in particular in the situation of our main Example 3.6, this minimal OOM is equivalent to an HMM and thus the minimal HMM can be obtained in a constructive way.

2 Predictive models of stochastic processes

2.1 Generating a process

Before suggesting our notion of prediction, we first consider the task of generating a process. Generating a predicted future based on memory states is a crucial part of our understanding of prediction. We assume a finite set D (called *state space* or *alphabet* of the generated process), a countable set M of *memory states* (also called *internal states*) and a Markov kernel (transition probability) which we call *generator*:

$$\text{gen} : M \rightarrow \mathcal{P}(D \times M),$$

⁴PSRs exist also for controlled systems, thus including actions of the observer. In principle, they can be non-linear, but most of the theory considers the linear case.

⁵The vector space consists of the signed measures on the internal states.

where $\mathcal{P}(\mathbb{D} \times \mathbb{M})$ is the set of probability measures on $\mathbb{D} \times \mathbb{M}$. The situation for more general spaces is discussed in the appendix.

Remark. A generator together with an initial probability distribution on the internal states is an (edge-emitting) HMM, except that we allow the set of internal states to be countable (instead of finite).

We use the notation $\text{gen}(m; x, \hat{m})$ to denote the probability of the pair (x, \hat{m}) with respect to $\text{gen}(m)$. Together with an initial probability distribution μ on the memory states \mathbb{M} , this kernel generates a stochastic process \tilde{X}_k , $k \in \mathbb{N}$, on \mathbb{D} and a process M_k , $k \in \mathbb{N}_0 := \mathbb{N} \cup \{0\}$, on \mathbb{M} in the following way: Being in a memory state at time k , it (stochastically) produces a new memory state at time $k + 1$ and, at the same time, emits a symbol from \mathbb{D} . This is shown in Figure 4.

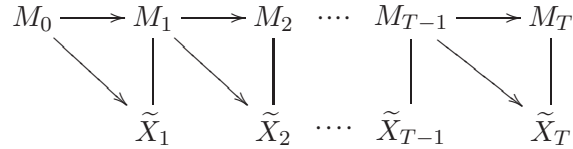


Figure 4: The process of generating memory states M_k and emitting observable states \tilde{X}_k .

The joint distribution is computed according to

$$\mathbb{P}(M_{[0,T]} = m_{[0,T]}, \tilde{X}_{[1,T]} = x_{[1,T]}) = \mu(m_0) \prod_{k=1}^T \text{gen}(m_{k-1}; x_k, m_k),$$

where we use the notation $[0, T]$ for the discrete interval $\{0, \dots, T\}$ and $M_{[0,T]} = m_{[0,T]}$ for $M_0 = m_0, \dots, M_T = m_T$. Similarly, throughout the paper we also use the notation $X_{\mathbb{T}}$ to denote a stochastic process X_k , $k \in \mathbb{T}$, where \mathbb{T} is the time set of the process.

Definition 2.1 (generating a process). Let $X_{\mathbb{N}}$ be a stochastic process on \mathbb{D} . We say that gen *generates* $X_{\mathbb{N}}$ if there exists an initial distribution μ for gen , such that $\tilde{X}_{\mathbb{N}}$ has the same distribution as $X_{\mathbb{N}}$.

Remark. For *every* stochastic process $X_{\mathbb{N}}$, there exists a generator gen which generates it. This is true because the set of internal states is allowed to be countable and the time set is only semi-infinite. Thus the generator can store the complete history of output symbols in the internal state. Another generator which can be constructed for any process is the finite-history version of the ε -machine (see Example 2.4).

2.2 Our prediction setting

We use generators as models for the process of prediction. The initial distribution is computed by a memory map from past observations and contains the information of the history. Although computational mechanics usually works with infinite length observations (histories), we allow here only finite but varying length observations. This way we avoid measure-theoretic

technicalities due to an uncountable set of infinite history trajectories. In the appendix, we treat the case of infinite histories and show that virtually everything remains valid. Unfortunately, the variation of history length leads to some notational technicalities, especially in Section 3.1.

Throughout this article, we consider a *stationary* stochastic process $X_{\mathbb{Z}}$, the *observable process*, with *finite* state set D . Note that, since $X_{\mathbb{Z}}$ is stationary, it is uniquely determined by its restriction to positive times. For the task of prediction, we assume that the outcome of $X_{\mathbb{Z}}$ is known for some finite but arbitrary past time interval $[-t+1, 0]$. Based on these observations, a generator is used as a mechanism for generating an outcome of $\tilde{X}_{[1,T]}^t$ as prediction of the real future outcome $X_{[1,T]}$. The situation is illustrated in Figure 5 and made more precise by the following definitions.

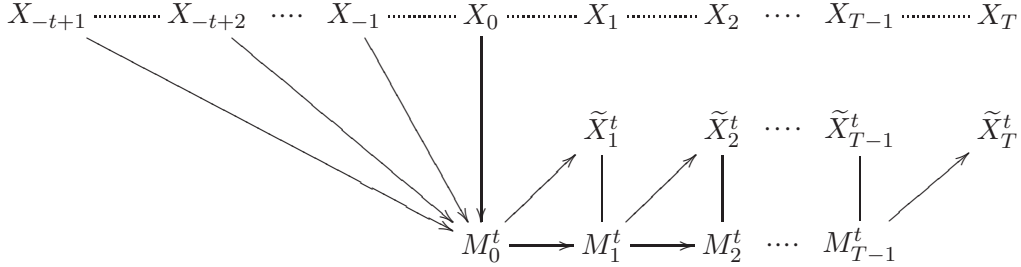


Figure 5: The process of generating $\tilde{X}_{[1,T]}^t$ as prediction of $X_{[1,T]}$ based on a length- t history which is fed into the memory variable M_0^t . The dotted lines symbolize that $X_{\mathbb{Z}}$ may have arbitrary dependencies and need not be Markov.

Definition 2.2 (memory). A *memory (map) mem* assigns to every history $x_{[-t+1,0]} \in D^{[-t+1,0]}$ of arbitrary but finite length t a probability distribution on a countable set M of *memory states*:

$$\text{mem}: D^* := \bigcup_{t \in \mathbb{N}_0} D^{[-t+1,0]} \rightarrow \mathcal{P}(M).$$

Note that D^* contains the “empty history,” which corresponds to not having observed anything.

We use the memory map mem and a generator $\text{gen}: M \rightarrow \mathcal{P}(D \times M)$ to define random variables M_k^t and \tilde{X}_k^t as shown in Figure 5. For every history length t , $X_{[-t+1,0]}$ and mem induce a random variable $M^t = M_0^t$ with distribution

$$\mathbb{P}(M^t = m \mid X_{[-t+1,0]} = x_{[-t+1,0]}) = \text{mem}(x_{[-t+1,0]}; m).$$

Now we can start the generator gen in the memory state M_0^t and obtain the predicted process $\tilde{X}_{\mathbb{N}}^t$ on D as well as a process of internal states $M_{\mathbb{N}_0}^t$ on M with the joint (conditional) distribution

$$\begin{aligned} \mathbb{P}(M_{[0,T]}^t = m_{[0,T]}, \tilde{X}_{[1,T]}^t = x_{[1,T]} \mid X_{[-t+1,0]} = x_{[-t+1,0]}) \\ = \text{mem}(x_{[-t+1,0]}; m_0) \prod_{k=1}^T \text{gen}(m_{k-1}; x_k, m_k). \end{aligned}$$

Definition 2.3 (predictive model). We call the memory map *mem predictive* (w.r.t. $X_{\mathbb{Z}}$) if there exists a generator $\text{gen}: \mathbb{M} \rightarrow \mathcal{P}(\mathbb{D} \times \mathbb{M})$, such that for all t and all $x_{[-t+1,0]}$ satisfying $\mathbb{P}(X_{[-t+1,0]} = x_{[-t+1,0]}) > 0$ the following equality holds:

$$\mathbb{P}(X_{[1,T]} \mid X_{[-t+1,0]} = x_{[-t+1,0]}) = \mathbb{P}(\tilde{X}_{[1,T]}^t \mid X_{[-t+1,0]} = x_{[-t+1,0]}) \quad \text{for all } T.$$

We then call the pair (mem, gen) *predictive model* of the process $X_{\mathbb{Z}}$.

This definition of a predictive memory corresponds to the requirement (7) which we already discussed in the introduction. To summarize, if we have a predictive model and a finite interval of observations with arbitrary length t , we use the memory map to (stochastically) produce an initial value M^t for the generator. Then we apply the generator to produce a predicted future $\tilde{X}_{[1,T]}^t$ that has the same statistic properties as the “real” future $X_{[1,T]}$, conditioned on the observations $X_{[-t+1,0]}$. It is important that the generator must not depend on the length t of the history.

In Section 3, we relate our notion of a predictive model to the definition used in computational mechanics. Before doing so, we give an example showing that one can construct a predictive model of *any* stationary stochastic process, namely (the finite-history version of) the ε -machine of computational mechanics. This important example is also used in Proposition 3.5 and Example 3.6.

Example 2.4 (ε -machine). In computational mechanics, the ε -*machine* is defined on the so-called *causal states*. These are defined as equivalence classes of observed histories. Usually these histories are assumed to have infinite length, but finite length histories have also been considered (e.g. [Feldman and Crutchfield, 1998]). In this case, the identified histories may have different lengths. The equivalence relation identifies histories with the same conditional expectation on the future, i.e. $x_{[-t+1,0]} \sim x'_{[-s+1,0]}$ if and only if⁶

$$\mathbb{P}(X_{[1,T]} \mid X_{[-t+1,0]} = x_{[-t+1,0]}) = \mathbb{P}(X_{[1,T]} \mid X_{[-s+1,0]} = x'_{[-s+1,0]}) \quad \text{for all } T > 0.$$

The *causal state* of $x_{[-t+1,0]}$ is given by its equivalence class

$$\mathfrak{C}(x_{[-t+1,0]}) := \{x'_{[-s+1,0]} \mid s \in \mathbb{N}_0, x'_{[-s+1,0]} \sim x_{[-t+1,0]}\}.$$

As memory set, we take the set $\mathbb{M}_{\mathfrak{C}} := \text{Im}(\mathfrak{C}) := \{\mathfrak{C}(x_{[-t+1,0]})\}$ of causal states. The deterministic memory map $\text{mem}_{\mathfrak{C}}$ assigns to a history $x_{[-t+1,0]}$ the Dirac measure in the corresponding causal state $\mathfrak{C}(x_{[-t+1,0]})$, i.e. $\text{mem}_{\mathfrak{C}}(x_{[-t+1,0]}; \mathfrak{C}(x_{[-t+1,0]})) = 1$. To get a predictive model, we also need a generator. By $x_{[-t+1,0]}x$, we denote the history $y_{[-t,0]}$ of length $t+1$ obtained by concatenation of $x_{[-t+1,0]}$ and x , i.e. $y_0 = x$ and $y_{-k} = x_{-k+1}$ for $k = 1, \dots, t$. Note that if $\mathfrak{C}(x_{[-t+1,0]}) = \mathfrak{C}(x'_{[-s+1,0]})$, we also have $\mathfrak{C}(x_{[-t+1,0]}x) = \mathfrak{C}(x'_{[-s+1,0]}x)$, provided that $x_{[-t+1,0]}$ and $x'_{[-s+1,0]}$ have positive probability. This is true because

$$\mathbb{P}(X_{[1,T]} \mid X_{[-t,0]} = x_{[-t+1,0]}x) = \frac{\mathbb{P}(X_0 = x, X_{[1,T]} \mid X_{[-t,-1]} = x_{[-t+1,0]})}{\mathbb{P}(X_0 = x \mid X_{[-t,-1]} = x_{[-t+1,0]})},$$

and $X_{\mathbb{Z}}$ is stationary. Therefore, the following generator (the ε -machine transition) is well defined:

$$\text{gen}_{\mathfrak{C}}(m; x, m') := \begin{cases} \mathbb{P}(X_1 = x \mid X_{[-t+1,0]} = x_{[-t+1,0]}), & \text{if } \mathfrak{C}(x_{[-t+1,0]}x) = m' \\ 0, & \text{otherwise} \end{cases},$$

⁶We assign histories with probability zero, i.e. $\mathbb{P}(X_{[-t+1,0]} = x_{[-t+1,0]}) = 0$, to arbitrary equivalence classes.

where $x_{[-t+1,0]}$ is any history with positive probability and $\mathfrak{C}(x_{[-t+1,0]}) = m$. We obtain

$$\begin{aligned}
& \mathbb{P}(\tilde{X}_{[1,T]}^t = x_{[1,T]} \mid X_{[-t+1,0]} = x_{[-t+1,0]}) \\
&= \prod_{k=1}^T \text{gen}_{\mathfrak{C}}(\mathfrak{C}(x_{[-t+1,0]}x_1 \cdots x_{k-1}); x_k, \mathfrak{C}(x_{[-t+1,0]}x_1 \cdots x_k)) \\
&= \prod_{k=1}^T \mathbb{P}(X_1 = x_k \mid X_{[-t-k+2,0]} = x_{[-t+1,0]}x_1 \cdots x_{k-1}) \\
&\stackrel{\text{(stationary)}}{=} \mathbb{P}(X_{[1,T]} = x_{[1,T]} \mid X_{[-t+1,0]} = x_{[-t+1,0]}).
\end{aligned}$$

Thus $(\text{mem}_{\mathfrak{C}}, \text{gen}_{\mathfrak{C}})$ is a predictive model. \diamond

As a pair (mem, gen) , a predictive model (of $X_{\mathbb{Z}}$) in particular provides the generator gen which generates the restriction $X_{\mathbb{N}}$ of the process $X_{\mathbb{Z}}$ to positive times in the sense of Definition 2.1. The appropriate initial distribution is given by the memory map for $t = 0$, i.e. by $\text{mem}(\emptyset)$, where \emptyset is the empty history. In the following proposition, we show the converse of this statement: Every generator that generates the positive time restriction $X_{\mathbb{N}}$ can be used in a predictive model with an appropriate memory map.⁷ In particular, if the number of memory states in \mathbb{M} is large enough to allow for generating the positive time restriction of the process, it is also large enough to admit a predictive model of $X_{\mathbb{Z}}$.

Proposition 2.5 (generator as predictive model). *Let $\text{gen}: \mathbb{M} \rightarrow \mathcal{P}(\mathbb{D} \times \mathbb{M})$ be a generator that generates the positive time restriction of the process $X_{\mathbb{Z}}$. Then there is a memory map $\text{mem}: \mathbb{D}^* \rightarrow \mathcal{P}(\mathbb{M})$, such that (mem, gen) is a predictive model of $X_{\mathbb{Z}}$.*

Proof. Let the initial distribution for gen be such that $\tilde{X}_{\mathbb{N}}$ has the same distribution as $X_{\mathbb{N}}$. Define for all $x_{[-t+1,0]}$ with positive probability

$$\text{mem}(x_{[-t+1,0]}; m) := \mathbb{P}(M_t = m \mid \tilde{X}_{[1,t]} = x_{[-t+1,0]}).$$

The case $t = 0$ is clear because of the fact that $\tilde{X}_{\mathbb{N}}^0$ and $\tilde{X}_{\mathbb{N}}$ have the same distribution. Therefore, let $t > 0$:

$$\begin{aligned}
& \mathbb{P}(\tilde{X}_{[1,T]}^t \mid X_{[-t+1,0]} = x_{[-t+1,0]}) \\
&= \sum_m \mathbb{P}(M_0^t = m \mid X_{[-t+1,0]} = x_{[-t+1,0]}) \mathbb{P}(\tilde{X}_{[1,T]}^t \mid M_0^t = m, X_{[-t+1,0]} = x_{[-t+1,0]}) \\
&= \sum_m \text{mem}(x_{[-t+1,0]}; m) \mathbb{P}(\tilde{X}_{[1,T]}^t \mid M_0^t = m) \\
&= \sum_m \mathbb{P}(M_t = m \mid \tilde{X}_{[1,t]} = x_{[-t+1,0]}) \mathbb{P}(\tilde{X}_{[t+1,t+T]} \mid M_t = m) \\
&= \mathbb{P}(\tilde{X}_{[t+1,t+T]} \mid \tilde{X}_{[1,t]} = x_{[-t+1,0]}) \stackrel{\text{(assumption)}}{=} \mathbb{P}(X_{[t+1,t+T]} \mid X_{[1,t]} = x_{[-t+1,0]}) \\
&\stackrel{\text{(stationary)}}{=} \mathbb{P}(X_{[1,T]} \mid X_{[-t+1,0]} = x_{[-t+1,0]}),
\end{aligned}$$

⁷Thus, generative and predictive models are essentially equivalent. Note, however, that the memory map need not be unique.

where we used that $\tilde{X}_{[1,T]}^t$ is independent of $X_{[-t+1,0]}$ given M_0^t and $\tilde{X}_{[t+1,t+T]}$ is independent of $\tilde{X}_{[1,t]}$ given M_t . \square

Remark. Proposition 2.5 is the main source of examples of predictive models: We can take an arbitrary (edge-emitting) HMM with stationary output process and automatically obtain a predictive model of the output. This model uses the same memory states and generating mechanism as the original HMM. Example 3.6 below is of that type.

3 What does “prediction” really mean?

3.1 Predictive versus prescient memories

As we already mentioned in the introduction, our concept of a predictive memory map differs from the concept usually discussed within computational mechanics. There, one tries to compress the observed sequence $x_{[-t+1,0]}$ by the memory map and requires that, at the same time, no information about the future $x_{[1,T]}$ (for all T) that is contained in the history $x_{[-t+1,0]}$ is lost. In the situation where all observed histories have the same length, which in computational mechanics is usually assumed to be infinite, this means requiring that the mutual information between history and future is equal to the mutual information between the memory variable M and the future, similar to the requirement (2) of the introduction. In our present setting of finite varying history lengths, however, we do not have a single memory state at time zero but for any history length t a different memory state M^t .⁸ Simply assuming the information equality for every length t separately, i.e.

$$I(M^t : X_{[1,T]}) = I(X_{[-t+1,0]} : X_{[1,T]}) \quad \text{for all } T \text{ and all } t, \quad (9)$$

is a weak requirement which does not provide the correct correspondence to (2) in the context of computational mechanics for finite but varying observation lengths. For memory maps satisfying (9), the information about the future need not be contained in the memory state alone but also in the particular observation length t . The same memory state m can have a completely different implication on the future if it results from different history lengths (see Example 3.1). Therefore, we have to assume that the memory keeps all information about the future without the additional knowledge of t . We give two equivalent versions of the right correspondence to (2).

First, we simply assume, in addition to (9), that conditional probabilities of the future given a memory state do not depend on the observation length t . More precisely, given $m \in \mathbb{M}$, we assume that $\mathbb{P}(X_{[1,T]} | M^t = m)$ is independent of t whenever $\mathbb{P}(M^t = m) > 0$. Since (9) is equivalent to

$$\mathbb{P}(X_{[1,T]} | X_{[-t+1,0]} = x_{[-t+1,0]}) = \mathbb{P}(X_{[1,T]} | M^t = m)$$

whenever $\mathbb{P}(X_{[-t+1,0]} = x_{[-t+1,0]}, M^t = m) > 0$, we finally get the following condition as correspondence to (2):

$$\begin{aligned} \mathbb{P}(X_{[-t+1,0]} = x_{[-t+1,0]}, M^t = m) > 0, \quad \mathbb{P}(M^s = m) > 0 \\ \implies \mathbb{P}(X_{[1,T]} | X_{[-t+1,0]} = x_{[-t+1,0]}) = \mathbb{P}(X_{[1,T]} | M^s = m) \quad \text{for all } T. \end{aligned} \quad (10)$$

⁸This technicality does not appear in the infinite-history case discussed in the appendix.

As a second definition, which is equivalent to (10), we provide a t -independent version of (9). To this end, we imagine that t is determined randomly by a \mathbb{N}_0 -valued random variable τ which is assumed to be independent of all other variables. We call such a variable τ *random time*. Combining the family of memory variables M^t , $t \in \mathbb{N}_0$, with a random time τ we get a new variable M^τ which is equal to M^t precisely when $\tau = t$. We require that, for all random times τ , the corresponding M^τ contains maximal information about the future, even if we don't know the value of τ . More precisely,

$$I(M^\tau : X_{[1,T]}) = I(X_{[-\tau+1,0]} : X_{[1,T]}) \quad \text{for all } T \in \mathbb{N} \text{ and all random times } \tau. \quad (11)$$

Note that (11) contains (9) as the special case of constant random times. It is straightforward to show (but omitted here) that (11) is equivalent to (10). We illustrate the difference between (9) and (11) by the following example.

Example 3.1 (the difference via random times). Let $X_{\mathbb{Z}}$ be a non-i.i.d. Markov process on $D := \{0, 1\}$. Define

$$M^t := X_0 \quad \text{and} \quad \widehat{M}^t := \begin{cases} X_0, & \text{if } t \text{ odd} \\ 1 - X_0, & \text{if } t \text{ even} \end{cases}.$$

Then both M and \widehat{M} satisfy (9), whereas M also satisfies (11) and \widehat{M} does not. This is because the information $\widehat{M}^\tau = m$ is useless if we don't know whether τ is odd or even. \diamond

Definition 3.2 (prescient). We call a memory *prescient*, if it satisfies the equivalent conditions (11) and (10).

The following example shows that predictive memories need not be prescient and can have fewer memory states than the ε -machine (which turns out to be the minimal prescient memory in Proposition 3.5). It is an extension of Example 1.1 to the setting of stochastic processes.

Example 3.3. Let $D := \{0, 1, 2\}$ and $X_{\mathbb{Z}}$ the Markov process defined by

$$\mathbb{P}(X_0 = x_0) = \frac{1}{3} \quad \text{and} \quad \mathbb{P}(X_{k+1} = x_{k+1} \mid X_k = x_k) = \begin{cases} \frac{2}{3}, & \text{if } x_{k+1} = x_k \neq 2 \\ \frac{1}{3}, & \text{if } x_k = 2 \text{ or } x_{k+1} = 2 \\ 0, & \text{otherwise} \end{cases}.$$

As the process is Markov and different last symbols of the history yield different expectations of the next symbol, there are three causal states as defined in Example 2.4 (the empty history induces the same expectation as the histories ending in the symbol 2). On the other hand, we obtain a predictive memory map with two memory states in the same way as in Example 1.1: Let $M := \{0, 1\}$ and

$$\text{mem}(x_{[-t+1,0]}; m) := \begin{cases} 1, & \text{if } t \neq 0 \text{ and } m = x_0 \neq 2 \\ \frac{1}{2}, & \text{if } t = 0 \text{ or } x_0 = 2 \\ 0, & \text{otherwise} \end{cases}.$$

Just as in the situation of Example 1.1, we see that mem is indeed predictive. \diamond

The property “prescient” only refers to the map \mathbf{mem} and not to the operational aspects of prediction given by the generating “mechanism” of a generator \mathbf{gen} . Therefore, it does not refer to the predicted process $\tilde{X}_{[1,T]}^t$. In the following proposition, we show that nevertheless we can associate to any prescient memory a generative mechanism: Every prescient memory is predictive. Furthermore, in the case of deterministic memory, predictive and prescient are equivalent.

Proposition 3.4 (predictive versus prescient).

1. *Every prescient memory map is predictive.*
2. *If a memory map is deterministic and predictive, then it is also prescient.*

Proof.

1. Assume w.l.o.g. that for all $m \in \mathbf{M}$ there is some t_m with $\mathbb{P}(M^{t_m} = m) > 0$ (otherwise, m may be removed from \mathbf{M}). Let \widehat{M}^t be constructed from $X_{[-t+2,1]}$ with \mathbf{mem} , just like M^t is constructed from $X_{[-t+1,0]}$. We define the generator

$$\mathbf{gen}(m; x, \hat{m}) := \mathbb{P}(\widehat{M}^{t_m+1} = \hat{m}, X_1 = x \mid M^{t_m} = m).$$

In view of Proposition 2.5, it suffices to show

$$\mathbb{P}(\tilde{X}_{[1,T]}^t = x_{[1,T]}) = \mathbb{P}(X_{[1,T]} = x_{[1,T]}).$$

We show the more general equation (for m with $\mathbb{P}(M_0^t = m) > 0$)

$$\mathbb{P}(\tilde{X}_{[1,T]}^t = x_{[1,T]} \mid M_0^t = m) = \mathbb{P}(X_{[1,T]} = x_{[1,T]} \mid M_0^t = m)$$

by induction over T . The case $T = 0$ is trivial. For $T > 0$:

$$\begin{aligned} & \mathbb{P}(\tilde{X}_{[1,T]}^t = x_{[1,T]} \mid M_0^t = m) & (12) \\ &= \sum_{\hat{m}} \mathbb{P}(M_1^t = \hat{m}, \tilde{X}_1^t = x_1 \mid M_0^t = m) \mathbb{P}(\tilde{X}_{[2,T]}^t = x_{[2,T]} \mid M_1^t = \hat{m}) \\ &= \sum_{\hat{m}} \mathbf{gen}(m; x_1, \hat{m}) \mathbb{P}(\tilde{X}_{[1,T-1]}^{t_m+1} = x_{[2,T]} \mid M_0^{t_m+1} = \hat{m}) \\ &\stackrel{(\text{ind. as.})}{=} \sum_{\hat{m}} \mathbb{P}(\widehat{M}^{t_m+1} = \hat{m}, X_1 = x_1 \mid M_0^{t_m} = m) \mathbb{P}(X_{[1,T-1]} = x_{[2,T]} \mid M_0^{t_m+1} = \hat{m}). \end{aligned}$$

Now using stationarity of $X_{\mathbb{Z}}$ and (10), which holds also for \widehat{M} instead of M due to stationarity, we obtain for those \hat{m} with $\mathbb{P}(\widehat{M}^{t_m+1} = \hat{m}, X_1 = x_1 \mid M_0^{t_m} = m) > 0$ that

$$\begin{aligned} \mathbb{P}(X_{[1,T-1]} = x_{[2,T]} \mid M^{t_m+1} = \hat{m}) &= \mathbb{P}(X_{[2,T]} = x_{[2,T]} \mid \widehat{M}^{t_m+1} = \hat{m}) & (13) \\ &= \mathbb{P}(X_{[2,T]} = x_{[2,T]} \mid \widehat{M}^{t_m+1} = \hat{m}, M_0^{t_m} = m, X_1 = x_1). \end{aligned}$$

In total we obtain the required equality:

$$\mathbb{P}(\tilde{X}_{[1,T]}^t \mid M_0^t = m) \stackrel{(12)+(13)}{=} \mathbb{P}(X_{[1,T]} \mid M_0^{t_m} = m) \stackrel{(10)}{=} \mathbb{P}(X_{[1,T]} \mid M_0^t = m).$$

2. Let $M^t = f_t(X_{[-t+1,0]})$. For $m = f(x_{[-t+1,0]}) = f(x'_{[-s+1,0]})$, we get from predictiveness

$$\begin{aligned} \mathbb{P}(X_{[1,T]} | X_{[-s+1,0]} = x'_{[-s+1,0]}) &= \mathbb{P}(\tilde{X}_{[1,T]}^s | X_{[-s+1,0]} = x'_{[-s+1,0]}) \\ &= \mathbb{P}(\tilde{X}_{[1,T]}^t | M^t = m) \\ &= \mathbb{P}(X_{[1,T]} | X_{[-t+1,0]} = x_{[-t+1,0]}) . \end{aligned}$$

Consequently, if $A := f_s^{-1}(m)$ is the set of length- s histories mapped to m ,

$$\begin{aligned} \mathbb{P}(X_{[1,T]} | M^s = m) &= \frac{\sum_{x'_{[-s+1,0]} \in A} \mathbb{P}(X_{[-s+1,0]} = x'_{[-s+1,0]}) \mathbb{P}(X_{[1,T]} | X_{[-s+1,0]} = x'_{[-s+1,0]})}{\sum_{x'_{[-s+1,0]} \in A} \mathbb{P}(X_{[-s+1,0]} = x'_{[-s+1,0]})} \\ &= \frac{\mathbb{P}(X_{[1,T]} | X_{[-t+1,0]} = x_{[-t+1,0]}) \sum \mathbb{P}(X_{[-s+1,0]} = x'_{[-s+1,0]})}{\sum \mathbb{P}(X_{[-s+1,0]} = x'_{[-s+1,0]})} \\ &= \mathbb{P}(X_{[1,T]} | X_{[-t+1,0]} = x_{[-t+1,0]}) . \end{aligned}$$

This is nothing but equation (10). □

Remark. Let mem be a predictive memory map. Consider the set $M' := \{\text{mem}(x) \mid x \in D^*\} \subseteq \mathcal{P}(M)$ of *information states* and the deterministic M' -valued memory map defined by $\text{mem}'(x; \text{mem}(x)) = 1$. Then mem' is easily seen to be a prescient memory map (i.e. a sufficient statistic on the past for the future). Note that while mem' is directly linked to mem , M' may be much bigger than M . In Proposition 3.5 we show that the cardinality of M' is bounded below by the number of causal states.

3.2 Implications on minimality of predictive models

Figure 6 illustrates the situation in view of Proposition 3.4 with the abbreviations “DM = deterministic memory,” “PsM = prescient memory,” and “PdM = predictive memory.” In

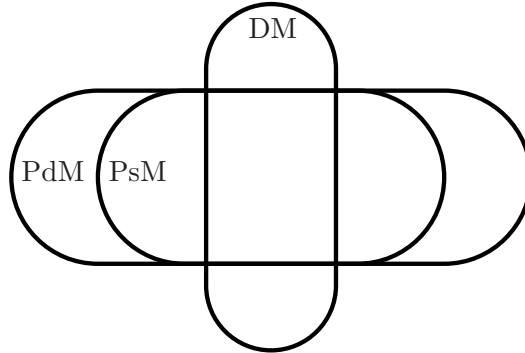


Figure 6: The extension of the memory class suggested by predictive models

computational mechanics, prescient deterministic memories have been studied, that is the intersection of DM and PsM. An extension of this intersection to larger classes of memory maps is natural. According to Proposition 3.4, we have the following hierarchy of possible extensions:

$$\text{DM} \cap \text{PsM} = \text{DM} \cap \text{PdM} \subsetneq \text{PsM} \subsetneq \text{PdM}. \quad (14)$$

According to the equality in (14), considering predictive memory maps without dropping the determinism requirement does not enlarge the class. Only recently, an extension to the class PsM including also non-deterministic memories has been considered by Still and Crutchfield [Still and Crutchfield, 2007]. It turns out that this extension does not allow for “smaller” models than already captured by deterministic memory maps, as we show in the following proposition. Therefore, we suggest to further extend the class from PsM to PdM and show in Example 3.6 that this extension is indeed effective.

Proposition 3.5 (ε -machine minimality in PsM). *The causal state projection \mathfrak{C} of Example 2.4 defines a prescient deterministic memory map to $\mathsf{M}_{\mathfrak{C}}$. Further, it has minimal number of memory states amongst the prescient (not necessarily deterministic) memories. More precisely:*

$$\text{mem}: \mathsf{D}^* \rightarrow \mathcal{P}(\mathsf{M}) \text{ prescient} \quad \Rightarrow \quad |\mathsf{M}| \geq |\mathsf{M}_{\mathfrak{C}}|.$$

Proof. We show (10). Let $m_{\mathfrak{C}} := \mathfrak{C}(x_{[-t+1,0]})$ and s be such that $\mathbb{P}(\mathfrak{C}(X_{[-s+1,0]}) = m_{\mathfrak{C}}) > 0$. The conditional probability $\mathbb{P}(X_{[1,T]} \mid \mathfrak{C}(X_{[-s+1,0]}) = m_{\mathfrak{C}})$ is a convex combination of $\mathbb{P}(X_{[1,T]} \mid X_{[-s+1,0]} = x'_{[-s+1,0]})$ with $x'_{[-s+1,0]} \in m_{\mathfrak{C}}$. Since all elements of $m_{\mathfrak{C}}$ induce the same conditional probability of $X_{[1,T]}$, we obtain

$$\mathbb{P}(X_{[1,T]} \mid X_{[-t+1,0]} = x_{[-t+1,0]}) = \mathbb{P}(X_{[1,T]} \mid \mathfrak{C}(X_{[-s+1,0]}) = m_{\mathfrak{C}})$$

and $\text{mem}_{\mathfrak{C}}$ is prescient. Now assume that mem is another prescient memory. We show that if the supports of $\text{mem}(x_{[-t+1,0]})$ and $\text{mem}(x'_{[-s+1,0]})$ are not disjoint, then $\mathfrak{C}(x_{[-t+1,0]}) = \mathfrak{C}(x'_{[-s+1,0]})$. In particular, $|\mathsf{M}_{\mathfrak{C}}| \leq |\mathsf{M}|$. Thus, assume some $m \in \mathsf{M}$ with

$$\mathbb{P}(M^t = m \mid X_{[-t+1,0]} = x_{[-t+1,0]}) > 0 \quad \text{and} \quad \mathbb{P}(M^s = m \mid X_{[-s+1,0]} = x'_{[-s+1,0]}) > 0.$$

From (10), we obtain

$$\mathbb{P}(X_{[1,T]} \mid X_{[-t+1,0]} = x_{[-t+1,0]}) = \mathbb{P}(X_{[1,T]} \mid M^s = m) = \mathbb{P}(X_{[1,T]} \mid X_{[-s+1,0]} = x'_{[-s+1,0]}),$$

hence $\mathfrak{C}(x_{[-t+1,0]}) = \mathfrak{C}(x'_{[-s+1,0]})$, which finishes the proof. \square

Example 3.6 below illustrates that our extension allows for minimal memories in PdM not captured within PsM. It is an example of an HMM with two hidden states that, consequently, admits a predictive model with two memory states. The corresponding minimal number of memory states within PsM, which according to Proposition 3.5 is realized by the ε -machine, turns out to be infinite.⁹ Moreover, the causal states are singletons, so that the causal state projection does not achieve any compression. Here, different histories lead to different expectations on the hidden states (different information states) of the HMM. When considering infinite histories, the causal states of the ε -machine turn out to be even uncountable, whereas the HMM can still be turned into a predictive model with two internal states (see the appendix).

⁹The existence of finite HMMs leading to infinitely many causal states was already mentioned in [Crutchfield, 1994].

Example 3.6 (predictive model smaller than the ε -machine). In order to specify this example, we consider an HMM (gen, μ) with internal state set M . This defines stochastic processes $\tilde{X}_{\mathbb{N}}$ and $M_{\mathbb{N}_0}$. If the joint process $(\tilde{X}_{\mathbb{N}}, M_{\mathbb{N}})$ is stationary, we can extend this joint process to a stationary process with time set \mathbb{Z} in a unique way. We denote the resulting processes on D and M with $X_{\mathbb{Z}}$ and $S_{\mathbb{Z}}$, respectively, and interpret them as the observable process and a process of internal states. In this concrete example, we take $D := M := \{0, 1\}$ and the uniform distribution on M as initial distribution. With a parameter p , $0 < p < \frac{1}{4}$, we define the generator by

$$\text{gen}(m; x, \hat{m}) := \begin{cases} 1 - 2p, & \text{if } \hat{m} = x = m \\ p, & \text{if } x \neq m \\ 0, & \text{otherwise} \end{cases}. \quad (15)$$

See Figure 7 for an illustration of the transition graph. It is easy to check that the stationarity

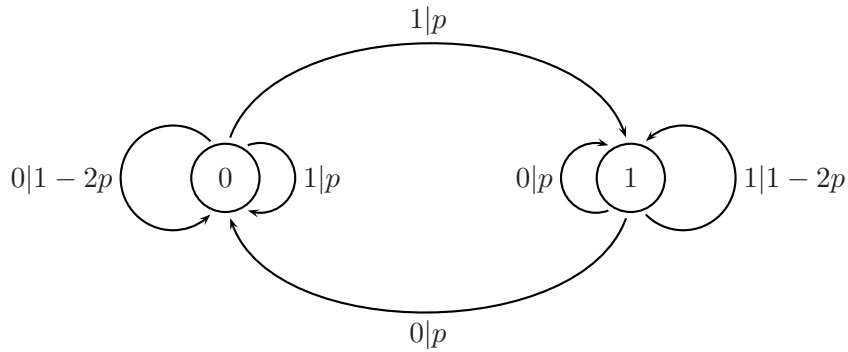


Figure 7: Transition graph of the generator defined by (15). Circled nodes are internal states and edges are transitions, labeled with output symbol x and transition probability q as “ $x|q$ ”.

condition is indeed satisfied. Because of Proposition 2.5, it is clear that there is a predictive model of the process $X_{\mathbb{Z}}$ with two memory states. We now show that, nevertheless, the causal states are singletons. For this purpose, we define for any output symbol $x \in D$ a function $f_x: [0, 1] \rightarrow [0, 1]$ which keeps track of the probability that the internal state is 0. Concretely,

$$f_x(y) := \frac{y \text{gen}(0; x, 0) + (1 - y) \text{gen}(1; x, 0)}{y \sum_{m=0}^1 \text{gen}(0; x, m) + (1 - y) \sum_{m=0}^1 \text{gen}(1; x, m)}.$$

We compute the conditional probability that the internal state is 0 as follows:

$$\begin{aligned} & \mathbb{P}(S_0 = 0 \mid X_{[-t+1,0]} = x_{[-t+1,0]}) \\ &= \sum_{m=0}^1 \frac{\mathbb{P}(S_{-1} = m \mid X_{[-t+1,-1]} = x_{[-t+1,-1]}) \mathbb{P}(S_0 = 0, X_0 = x_0 \mid S_{-1} = m)}{\mathbb{P}(X_0 = x_0 \mid X_{[-t+1,-1]} = x_{[-t+1,-1]})} \\ &= f_{x_0}(\mathbb{P}(S_{-1} = 0 \mid X_{[-t+1,-1]} = x_{[-t+1,-1]})) \\ &\stackrel{(\text{induction})}{=} f_{x_0} \circ \dots \circ f_{x_{-t+1}}(\mathbb{P}(S_{-t} = 0)) = f_{x_0} \circ \dots \circ f_{x_{-t+1}}(\tfrac{1}{2}) \end{aligned}$$

Obviously, $\mathbb{P}(X_{[1,T]} \mid S_0 = 0) \neq \mathbb{P}(X_{[1,T]} \mid S_0 = 1)$ (as $p \neq \frac{1}{4}$), and therefore

$$\begin{aligned} \mathbb{P}(X_{[1,T]} \mid X_{[-t+1,0]} = x_{[-t+1,0]}) &= \mathbb{P}(X_{[1,T]} \mid X_{[-s+1,0]} = x'_{[-s+1,0]}) \\ &\Leftrightarrow f_{x_0} \circ \dots \circ f_{x_{-t+1}}(\tfrac{1}{2}) = f_{x'_0} \circ \dots \circ f_{x'_{-s+1}}(\tfrac{1}{2}). \quad (16) \end{aligned}$$

Now plugging the definition of gen into the definition of f_x we obtain

$$f_0(y) = \frac{y(1-3p)+p}{y(1-4p)+2p} \quad \text{and} \quad f_1(y) = \frac{yp}{1-2p-y(1-4p)}.$$

We observe that both f_0 and f_1 are strictly increasing,

$$f_0(]0, 1]) =]\frac{1}{2}, 1[\quad \text{and} \quad f_1(]0, 1]) =]0, \frac{1}{2}[.$$

This implies that $f_{x_0} \circ \dots \circ f_{x_{-t+1}}(\frac{1}{2})$ and $f_{x'_0} \circ \dots \circ f_{x'_{-s+1}}(\frac{1}{2})$ are different for distinct $x_{[-t+1,0]}$ and $x'_{[-s+1,0]}$. Because of (16), the causal states are singletons. \diamond

4 Conclusions

There are two natural and well-known concepts of constructing an HMM as model of a given stochastic process. The first one is based on the idea that the memory variable (hidden variable of the HMM) is prescient in the sense of being a sufficient statistic on the past for the future. Minimizing the size of such a prescient memory leads to the so-called ε -machine which can be obtained in a constructive and efficient way. The second way of assigning an HMM to a stochastic process requires the HMM to be capable of generating the given process. The problem of finding a minimal such HMM is an intrinsically hard problem and can not be solved as easily as in the ε -machine case. Both approaches provide different ways of understanding prediction which is a central theme of computational mechanics. We feel that the current literature does not highlight the difference between these two approaches in sufficient detail. This, unfortunately, leaves space for a misperception of ε -machines as minimal HMMs capable of generating the underlying process, which corresponds to the second approach and not to the ε -machine approach. Based on this perception, there is an apparent inconsistency between the ε -machine minimality and the existence of (substantially) smaller hidden Markov models that generate the same underlying process. In this article, we clarified the formal relation between the two approaches and thereby identified the ε -machine approach as the more restrictive one. We argued that the assumptions made within the second approach are related to a generative understanding of prediction which leads to the notion of a predictive memory and a corresponding predictive model. Based on this generative understanding of prediction we could show that the underlying HMM can be significantly smaller than the corresponding ε -machine. We compared the operational aspects of the two concepts of prediction and discussed some drawbacks of our approach in this regard. Currently, we do not know which of the two notions of prediction is the more natural one and further steps towards revealing and comparing operational aspects of prediction are subject of our research.

A Appendix: Prediction with infinite histories

In this appendix, we translate our setting and results to the case of infinite histories. Essentially all results remain true, but we have to deal with measure theoretic issues because the set of infinite histories is uncountable and may enforce uncountably many internal states of predictive models including the ε -machine. As it does not require much additional effort, we also consider more general, non-discrete state spaces for the (observable) stochastic process.

A.1 Generating a process

We assume that the state space D of the generated process is a Polish space, the space M of memory (or internal) states is a Souslin space,¹⁰ and the generator gen is a Markov kernel, i.e. a *measurable* map

$$\text{gen}: M \rightarrow \mathcal{P}(D \times M).$$

Here, the topological spaces M and D are equipped with their Borel σ -algebras and $\mathcal{P}(D \times M)$ with the σ -algebra induced by the evaluations $e_A: \mu \mapsto \mu(A)$ which coincides with the Borel σ -algebra of the weak-* topology. We use the notation $\text{gen}(m; D \times B)$ to denote the probability of a measurable set (event) $D \times B$ with respect to $\text{gen}(m)$. Like in the discrete case, gen together with an initial distribution μ on M generates stochastic processes $\tilde{X}_{\mathbb{N}}$ and $M_{\mathbb{N}_0}$ on D and M (the Kolmogorov extension theorem holds for Souslin spaces). The finite dimensional joint distributions are computed according to

$$\begin{aligned} & \mathbb{P}(\tilde{X}_{[1,T]} \in D_{[1,T]}, M_{[0,T]} \in B_{[0,T]}) \\ &= \int_{B_0} \int_{D_1 \times B_1} \cdots \int_{D_T \times B_T} 1 \text{gen}(m_{T-1}; d(d_T, m_T)) \cdots \text{gen}(m_0; d(d_1, m_1)) \mu(dm_0), \end{aligned}$$

where we use the notation $D_I = \prod_{i \in I} D_i$ for measurable sets D_i . If the initial distribution μ is gen -invariant, the processes $M_{\mathbb{N}_0}$ and $\tilde{X}_{\mathbb{N}}$ are extended to stationary processes $M_{\mathbb{Z}}$ and $\tilde{X}_{\mathbb{Z}}$.

Definition A.1 (generating a process). Let $X_{\mathbb{N}}$ (resp. $X_{\mathbb{Z}}$) be a stochastic process on D . We say that gen *generates* $X_{\mathbb{N}}$ (resp. $X_{\mathbb{Z}}$) if there exists an initial distribution (resp. *invariant* distribution) for gen such that $\tilde{X}_{\mathbb{N}}$ (resp. $\tilde{X}_{\mathbb{Z}}$) has the same distribution as $X_{\mathbb{N}}$ (resp. $X_{\mathbb{Z}}$).

A.2 Our prediction setting

Let $X_{\mathbb{Z}}$ be a D -valued, stationary stochastic process and D Polish.

Definition A.2 (memory). A *memory (map)* mem is a Markov kernel from the set $D^{-\mathbb{N}_0}$ of history trajectories to a Souslin space M of *memory states*:

$$\text{mem}: D^{-\mathbb{N}_0} \rightarrow \mathcal{P}(M).$$

As in Section 2.2, mem induces a random variable $M = M_0$ and, together with a generator gen , we obtain processes $\tilde{X}_{\mathbb{N}}$, $M_{\mathbb{N}_0}$. Note that here we have a single memory state at time 0 and there is no need for an upper index, as the history length is fixed (and infinite). In the following, we will use conditional probabilities and make the general assumption that they are *regular* versions. In Souslin spaces, this is always possible. Also note that the σ -algebra of a Souslin space is countably generated.

¹⁰Souslin spaces are slightly more general than Polish spaces ensuring that images of Souslin spaces under measurable maps, such as the set of causal states, are again Souslin spaces (see [Bourbaki, 1989]).

Definition A.3 (predictive model). We call the memory map mem *predictive* (w.r.t. $X_{\mathbb{Z}}$) if there exists a generator $\text{gen}: \mathbb{M} \rightarrow \mathcal{P}(\mathbb{D} \times \mathbb{M})$, such that

$$\mathbb{P}(X_{\mathbb{N}} | X_{-\mathbb{N}_0}) = \mathbb{P}(\tilde{X}_{\mathbb{N}} | X_{-\mathbb{N}_0}) \quad \text{a.s. (almost surely)}$$

We then call the pair (mem, gen) *predictive model* of the process $X_{\mathbb{Z}}$.

Proposition A.4 (generator as predictive model). *Let $\text{gen}: \mathbb{M} \rightarrow \mathcal{P}(\mathbb{D} \times \mathbb{M})$ be a generator that generates the process $X_{\mathbb{Z}}$. Then there is a memory map $\text{mem}: \mathbb{D}^{-\mathbb{N}_0} \rightarrow \mathcal{P}(\mathbb{M})$, such that (mem, gen) is a predictive model of $X_{\mathbb{Z}}$.*

Proof. Let the invariant initial distribution for gen be such that $\tilde{X}_{\mathbb{Z}}$ has the same distribution as $X_{\mathbb{Z}}$. Define mem by

$$\text{mem} \circ \tilde{X}_{-\mathbb{N}_0} := \mathbb{P}(M_0 | \tilde{X}_{-\mathbb{N}_0}).$$

We denote the processes generated by gen started in $\text{mem} \circ X_{-\mathbb{N}_0}$ by $\widehat{X}_{\mathbb{N}}$ and $\widehat{M}_{\mathbb{N}_0}$, respectively. We obtain

$$\begin{aligned} \mathbb{P}(\widehat{X}_{\mathbb{N}} | X_{-\mathbb{N}_0} = x) &= \int \mathbb{P}(\widehat{X}_{\mathbb{N}} | \widehat{M}_0 = m) \text{mem}(x; dm) \\ &= \int \mathbb{P}(\tilde{X}_{\mathbb{N}} | M_0 = m) \mathbb{P}(M_0 \in dm | \tilde{X}_{-\mathbb{N}_0} = x) \\ &= \mathbb{P}(\tilde{X}_{\mathbb{N}} | \tilde{X}_{-\mathbb{N}_0} = x) = \mathbb{P}(X_{\mathbb{N}} | X_{-\mathbb{N}_0} = x), \end{aligned}$$

$\mathbb{P} \circ X_{-\mathbb{N}_0}^{-1}$ -a.s. in $x \in \mathbb{D}^{-\mathbb{N}_0}$. Thus (mem, gen) is a predictive model. \square

Remark. Note that in Proposition A.4, in contrast to the finite-history case (see Proposition 2.5), it is not sufficient to assume that gen generates only the positive time process $X_{\mathbb{N}}$. Instead, we have to assume that the whole process $X_{\mathbb{Z}}$ is generated. This is, however, satisfied in our main Example 3.6. Thus it also admits a predictive model for infinite histories.

A.3 Comparison to computational mechanics

Definition A.5 (prescient). We call a memory map mem *prescient*, if the resulting M satisfies

$$\mathbb{P}(X_{\mathbb{N}} | M) = \mathbb{P}(X_{\mathbb{N}} | X_{-\mathbb{N}_0}) \quad \text{a.s.}$$

Predictive and prescient memories are related in the same way as in the finite-history case:

Proposition A.6 (predictive versus prescient).

1. *Every prescient memory map is predictive.*
2. *If a memory map is deterministic and predictive, then it is also prescient.*

Proof.

1. Let mem be a prescient memory and $M^k = M_0^k$ the resulting memory state at time k (calculated by mem from $X_{] -\infty, k]}$). Define the generator $\text{gen}: \mathbb{M} \rightarrow \mathcal{P}(\mathbb{D} \times \mathbb{M})$ by

$$\text{gen}(m; D \times B) := \mathbb{P}(X_1 \in D, M^1 \in B | M^0 = m), \quad (17)$$

and let $M_l = M_l^0$ be the memory state produced by gen at time l when started in M^0 at time 0. We show that $\mathbb{P}(\tilde{X}_{[1,T]} | M^0) = \mathbb{P}(X_{[1,T]} | M^0)$ a.s. by induction over T . We obtain

$$\begin{aligned} \mathbb{P}(\tilde{X}_{[2,T]} | M_1 = m) &= \mathbb{P}(\tilde{X}_{[1,T-1]} | M_0 = m) \stackrel{(\text{ind. as.})}{=} \mathbb{P}(X_{[1,T-1]} | M_0 = m) \\ &\stackrel{(X_{\mathbb{Z}} \text{ stat.})}{=} \mathbb{P}(X_{[2,T]} | M_1 = m), \end{aligned} \tag{18}$$

which leads to

$$\begin{aligned} \mathbb{P}(\tilde{X}_{[1,T]} \in D_{[1,T]} | M^0) &= \int_{D_1 \times M} \mathbb{P}(\tilde{X}_{[2,T]} \in D_{[2,T]} | M_1^0 = \cdot) \, d\text{gen}(M^0) \\ &\stackrel{(17)+(18)}{=} \int_{X_1^{-1}(D_1)} \mathbb{P}(X_{[2,T]} \in D_{[2,T]} | M^1) \, d\mathbb{P}(\cdot | M^0) \\ &\stackrel{(\text{prescience})}{=} \mathbb{P}(X_{[1,T]} \in D_{[1,T]} | M^0). \end{aligned}$$

This finishes the induction and in total we obtain

$$\mathbb{P}(\tilde{X}_{\mathbb{N}} | X_{-\mathbb{N}_0}) = \int \overbrace{\mathbb{P}(\tilde{X}_{\mathbb{N}} | M^0)}{=\mathbb{P}(X_{\mathbb{N}} | M^0)} \, d\mathbb{P}(\cdot | X_{-\mathbb{N}_0}) \stackrel{(\text{prescience})}{=} \mathbb{P}(X_{\mathbb{N}} | X_{-\mathbb{N}_0}).$$

2. Now let (mem, gen) be a predictive model and mem deterministic, i.e. $M = M^0 = f \circ X_{-\mathbb{N}_0}$ for some measurable function $f: \mathbf{D}^{-\mathbb{N}_0} \rightarrow \mathbf{M}$. Then $\mathbb{P}(X_{\mathbb{N}} | X_{-\mathbb{N}_0}) = \mathbb{P}(\tilde{X}_{\mathbb{N}} | M)$ is $\sigma(M)$ -measurable modulo \mathbb{P} . Since obviously $\sigma(M) \subseteq \sigma(X_{-\mathbb{N}_0})$, it holds a.s. that

$$\mathbb{P}(X_{\mathbb{N}} | X_{-\mathbb{N}_0}) = \mathbb{P}(X_{\mathbb{N}} | M). \quad \square$$

The causal states are defined as equivalence classes of histories inducing the same expectation on the future (with a fixed, regular version of conditional probability). Within the measure theoretic setting of this appendix, we find it more convenient to use a different but equivalent definition. To this end, consider the *causal state projection* $\mathfrak{C}: \mathbf{D}^{-\mathbb{N}_0} \rightarrow \mathcal{P}(\mathbf{D}^{\mathbb{N}})$ defined by

$$x \mapsto \mathfrak{C}(x) := \mathbb{P}(X_{\mathbb{N}} | X_{-\mathbb{N}_0} = x).$$

The set of *causal states* is the image of this map:

$$\mathbf{M}_{\mathfrak{C}} := \text{Im}(\mathfrak{C}) \subseteq \mathcal{P}(\mathbf{D}^{\mathbb{N}}).$$

Note that $\mathbf{M}_{\mathfrak{C}}$ inherits a Souslin topology and thus a measurable structure from $\mathcal{P}(\mathbf{D}^{\mathbb{N}})$. The map \mathfrak{C} is measurable, and it is easy to see that the corresponding deterministic memory is prescient. Before we prove that \mathfrak{C} corresponds to, as in the discrete case, a minimal prescient memory, we consider the process $X_{\mathbb{Z}}$ of Example 3.6. With a slight modification of the arguments there, one can show that the corresponding set $\mathbf{M}_{\mathfrak{C}}$ of causal states is uncountable (regardless of the version of conditional probability). On the other hand, there is a predictive model (for infinite histories) with two memory states.

Proposition A.7 (ε -machine minimality). *Let mem be a prescient memory map with set M of memory states. Then there exist disjoint measurable subsets $N_z \subseteq M$, $z \in M_{\mathfrak{C}}$, such that*

$$\text{mem}(x; N_{\mathfrak{C}(x)}) = 1 \quad \text{a.s.}$$

In particular, M cannot be essentially (i.e. up to zero-sets) smaller than $M_{\mathfrak{C}}$.

Proof. Let $M = M_0$ be the memory variable obtained by mem . Define $P_M: M \rightarrow \mathcal{P}(D^{\mathbb{N}})$ by

$$P_M \circ M := \mathbb{P}(X_{\mathbb{N}} | M), \quad \text{and} \quad N_z := P_M^{-1}(z) \text{ for } z \in M_{\mathfrak{C}}.$$

The N_z are obviously measurable and disjoint. Due to prescience, $P_M \circ M = \mathbb{P}(X_{\mathbb{N}} | X_{-\mathbb{N}_0})$ a.s., and consequently

$$\text{mem}(x; N_{\mathfrak{C}(x)}) = \mathbb{P}(M \in N_{\mathfrak{C}(x)} | X_{-\mathbb{N}_0} = x) = \mathbb{P}(P_M \circ M = \mathfrak{C}(x) | X_{-\mathbb{N}_0} = x) = 1,$$

where we used $\mathfrak{C}(x) = \mathbb{P}(X_{\mathbb{N}} | X_{-\mathbb{N}_0} = x)$. □

Acknowledgments

The authors are grateful for discussions with Jim Crutchfield, Cosma Shalizi, and Susanne Still. They also thank the anonymous referees for their constructive comments which substantially improved the paper. Nihat Ay thanks the Santa Fe Institute for hosting him during the initial work on this paper.

References

- [Ay and Crutchfield, 2005] Ay, N. and Crutchfield, J. P. (2005). Reductions of hidden information sources. *Journal of Statistical Physics*, 120:659–684.
- [Bialek et al., 2001] Bialek, W., Nemenman, I., and Tishby, N. (2001). Predictability, complexity, and learning. *Neural Computation*, 13:2409–2463.
- [Bourbaki, 1989] Bourbaki, N. (1989). *General Topology, Chapters 5-10*. Springer-Verlag.
- [Bukharaev, 1995] Bukharaev, R. G. (1995). *Theorie der stochastischen Automaten*. B.G. Teubner.
- [Crutchfield, 1994] Crutchfield, J. P. (1994). The calculi of emergence: Computation, dynamics and induction. *Physica D*, 75:11–54.
- [Crutchfield and Young, 1989] Crutchfield, J. P. and Young, K. (1989). Inferring statistical complexity. *Phys. Rev. Lett.*, 63:105–108.
- [Faigle and Schönhuth, 2007] Faigle, U. and Schönhuth, A. (2007). Asymptotic mean stationarity of sources with finite evolution dimension. *IEEE Transactions on Information Theory*, 53(7):2342–2348.

- [Feldman and Crutchfield, 1998] Feldman, D. P. and Crutchfield, J. P. (1998). Discovering noncritical organization: Statistical mechanical, information theoretic, and computational views of patterns in one-dimensional spin systems. Santa Fe Institute Working Paper 98-04-026.
- [Grassberger, 1986] Grassberger, P. (1986). Toward a quantitative theory of self-generated complexity. *Int. J. Theor. Phys.*, 25:907–938.
- [Heller, 1965] Heller, A. (1965). On stochastic processes derived from markov chains. *Annals of Mathematical Statistics*, 36:1286–1291.
- [Jaeger, 2000] Jaeger, H. (2000). Observable operator models for discrete stochastic time series. *Neural Computation*, 12(6):1371–1398.
- [Kulhavý, 1996] Kulhavý, R. (1996). *Recursive Nonlinear Estimation: A Geometric Approach*, volume 216 of *Lecture Notes in Control and Information Sciences*. Springer.
- [Littman et al., 2001] Littman, M. L., Sutton, R. S., and Singh, S. (2001). Predictive representations of state. *Advances in Neural Information Processing Systems*, 14.
- [Schönhuth and Jaeger, 2007] Schönhuth, A. and Jaeger, H. (2007). Characterization of ergodic hidden markov sources. Technical report, Zentrum für Angewandte Informatik Köln. Submitted to IEEE Transactions on Information Theory.
- [Shalizi and Crutchfield, 2001] Shalizi, C. R. and Crutchfield, J. P. (2001). Computational mechanics: Pattern and prediction, structure and simplicity. *Journal of Statistical Physics*, 104:817–879.
- [Shalizi et al., 2002] Shalizi, C. R., Shalizi, K. L., and Crutchfield, J. P. (2002). An algorithm for pattern discovery in time series. Informal publication, <http://arxiv.org/abs/cs.LG/0210025>.
- [Singh et al., 2004] Singh, S., James, M. R., and Rudary, M. R. (2004). Predictive state representations: a new theory for modeling dynamical systems. In *AUAI '04: Proceedings of the 20th conference on Uncertainty in artificial intelligence*, pages 512–519, Arlington, Virginia, United States. AUAI Press.
- [Still and Crutchfield, 2007] Still, S. and Crutchfield, J. P. (2007). Optimal causal inference. Informal publication, <http://arxiv.org/abs/0708.1580>.
- [Tishby et al., 1999] Tishby, N., Pereira, F. C., and Bialek, W. (1999). The information bottleneck method. Proceedings of the 37th Annual Allerton Conference on Communication, Control, and Computing, 358-377.