

Temporal Infomax Leads to Almost Deterministic Dynamical Systems

Nihat Ay and Thomas Wennekers¹

*Max-Planck-Institute for Mathematics in the Sciences,
Inselstr. 22-26, D-04103 Leipzig, Germany*

Abstract

The well-known Kullback-Leibler divergence of a random field from its factorization quantifies spatial interdependencies of the corresponding stochastic elements. In the present work, we introduce a generalized measure for stochastic interaction that captures also temporal interdependencies. Maximization of stochastic interaction in the setting of Markov chains is shown analytically and by simulations to result in an almost deterministic global dynamics, but, simultaneously, almost unpredictable single unit activity.

Key words: Markov model; Stochastic interaction; Information maximization;

1 Introduction

Information processing in the brain essentially depends on interrelations between neurons expressed by cooperation and competition in neural assemblies. The generation of spatio-temporal activation patterns, for instance, is at the base of statistical analyses as well as functional theories for neural synchronization and associative propagation of activity [1,6,9,11]. The hypothesis of strong interrelations is contained also in many conceptual approaches to an understanding of first principles for neural organization and learning. Information theory provides an appropriate framework for an analysis of such principles [4,7,8,10]. A well-known measure that quantifies relations of interacting units is the so-called *mutual information*: The Kullback-Leibler divergence

$$I(p) := D(p \| p_1 \otimes \cdots \otimes p_N) = \sum_{\nu=1}^N H(p_\nu) - H(p), \quad (1)$$

¹ Tel.: +49-341-9959-558; Fax: +49-341-9959-555, Email: nay@mis.mpg.de

where $H(\cdot)$ denotes the Shannon entropy and p_ν the ν 'th marginal of p , measures the “distance” of p from the set of factorized random fields. It is a natural measure for “spatial” interdependence and a starting point of many approaches to neural complexity. In [2,3] it has been studied from the *information geometric* point of view, where it is referred to as (*stochastic*) *interaction*.

In order to capture the intrinsically temporal aspects of dynamic interaction, I in (1) has been extended by Ay (cf. [4]) to the dynamical setting of Markov transitions. In the present paper we consider processes that optimize this temporal version of stochastic interaction. This leads to analytical results concerning the most fundamental feature of strongly interacting stochastic systems, i.e., the development of the system dynamics towards determinism, and provides a necessary first step for an investigation of learning processes also in more detailed recurrent neural network models.

2 Theoretical Framework and Results

Consider the set $V = \{1, \dots, N\}$ of binary units with the state sets $\Omega_\nu = \{0, 1\}$, $\nu \in V$. For a subsystem $A \subset V$, $\Omega_A := \{0, 1\}^A$ denotes the set of all configurations on A , and $\bar{P}(\Omega_A)$ is the set of probability distributions on Ω_A . Given two subsets A and B , where B is non-empty, $\bar{K}(\Omega_B | \Omega_A)$ is the set of all Markov transition kernels from Ω_A to Ω_B . In the case $A = B$, we use the abbreviation $\bar{K}(\Omega_A)$.

For a probability distribution $p \in \bar{P}(\Omega_A)$ and a Markov kernel $K \in \bar{K}(\Omega_B | \Omega_A)$, the *conditional entropy* of (p, K) is defined as

$$H(p, K) = - \sum_{\omega \in \Omega_A, \omega' \in \Omega_B} p(\omega) K(\omega' | \omega) \ln K(\omega' | \omega). \quad (2)$$

We set $H(Y | X) := H(p, K)$ if X and Y are random variables with $\text{Prob}\{X = \omega, Y = \omega'\} = p(\omega) K(\omega' | \omega)$ for all $\omega \in \Omega_A$ and $\omega' \in \Omega_B$. In the following we mainly consider the case $A = B = V$ where $H(p, K)$ measures the average uncertainty about the next state of the system given the present state. Assume p is strictly positive. Then H vanishes iff K is deterministic — that is, if the support set, $\text{supp } K(\cdot | \omega)$, contains only a single possible transition, $|\text{supp } K(\cdot | \omega)| = 1$, for all ω . H attains its maximal value $N \ln 2$ iff $K(\omega' | \omega) = 2^{-N}$ for all $\omega, \omega' \in \Omega_V$.

A Markov kernel $K \in \bar{K}(\Omega_V)$ is called *parallel* if there exist kernels $K^{(\nu)} \in \bar{K}(\Omega_\nu | \Omega_V)$, $\nu \in V$, such that

$$K(\omega' | \omega) = \prod_{\nu \in V} K^{(\nu)}(\omega'_\nu | \omega), \quad \text{for all } \omega, \omega' \in \Omega_V. \quad (3)$$

A kernel $K \in \bar{K}(\Omega_V)$ is called *split* if there exist kernels $K^{(\nu)} \in \bar{K}(\Omega_\nu)$, $\nu \in V$, such that

$$K(\omega' | \omega) = \prod_{\nu \in V} K^{(\nu)}(\omega'_\nu | \omega_\nu), \quad \text{for all } \omega, \omega' \in \Omega_V. \quad (4)$$

The split kernels are a proper subset of the parallel kernels and represent a dynamical version of the factorizable probability distributions. Thus, in analogy to (1) we define the *stochastic interaction* of the units with respect to a distribution $p \in \bar{P}(\Omega_V)$, p strictly positive, and a transition kernel $K \in \bar{K}(\Omega_V)$ as the p -divergence of K from being split: For this purpose, we define the marginal kernels $K_\nu \in K(\Omega_\nu)$, $\nu \in V$, of K by

$$K_\nu(\omega'_\nu | \omega_\nu) := \frac{\sum_{\substack{\sigma, \sigma' \in \Omega_V \\ \sigma_\nu = \omega_\nu, \sigma'_\nu = \omega'_\nu}} p(\sigma) K(\sigma' | \sigma)}{\sum_{\substack{\sigma \in \Omega_V \\ \sigma_\nu = \omega_\nu}} p(\sigma)}, \quad \omega_\nu, \omega'_\nu \in \Omega_\nu. \quad (5)$$

Then, the *stochastic interaction measure* I is given by the continuous extension of the function

$$I(p, K) := \sum_{\nu \in V} H(p_\nu, K_\nu) - H(p, K). \quad (6)$$

to the set $\bar{P}(\Omega_V) \times \bar{K}(\Omega_V)$ (cf. [4]). Equation (6) generalizes (1) to Markov transitions. $I(p, K)$ is large if the marginal transitions have high entropy, but that of the full transition is low. Supposed the current state $\omega \in \Omega_V$ is known, this corresponds with high predictability of the next state, but, conversely, not much information is gained from knowledge about single units, ω_ν . The predictability and, hence, degree of determinism of systems that optimize I is further characterized by the following Theorem.

Theorem 1: Consider a probability distribution $p \in \bar{P}(\Omega_V)$ and a transition kernel $K \in \bar{K}(\Omega_V)$. If (p, K) is a local maximizer of the interaction measure I then for all $\omega \in \text{supp } p$ one has $|\text{supp } K(\cdot | \omega)| \leq 1 + N$.

Theorem 1 presents a restriction to binary units of a more general Theorem proved in [5]. Note, that the expression $|\text{supp } K(\cdot | \omega)|$ counts the number of transitions with non-vanishing probability from an arbitrary state $\omega \in \text{supp } p$ to its target states. Since, there are exponentially many possible target states, namely 2^N , the estimate in Theorem 1, which is only linear in N , provides a strong upper bound on the number of transitions in strongly interacting systems. This has a direct implication to the determinism in such systems:

Corollary 2: In the Situation of Theorem 1, the entropy of (p, K) can be estimated by $H(p, K) \leq \ln(1 + N)$.

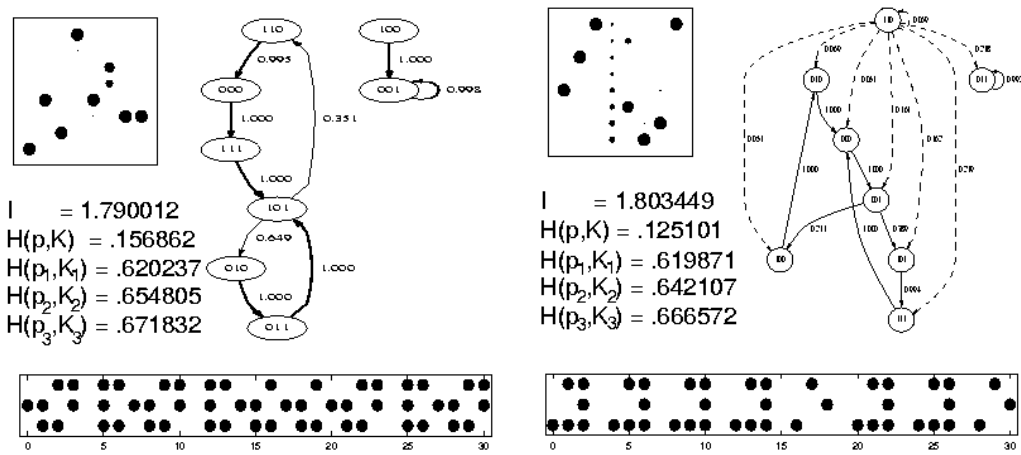


Fig. 1. Markov matrix (upper left), state transition graph (right), and sample activity (bottom) for two systems of each $N = 3$ strongly interacting units.

3 Simulations

Consider a Markov chain $X_n = (X_{\nu, n})_{\nu \in V}$, $n = 0, 1, 2, \dots$, given by an initial distribution $p \in \bar{P}(\Omega_V)$ and a kernel $K \in \bar{K}(\Omega_V)$. A probability distribution $p \in \bar{P}(\Omega_V)$ is called *stationary* with respect to $K \in \bar{K}(\Omega_V)$ iff $\sum_{\omega \in \Omega_V} p(\omega) K(\omega' | \omega) = p(\omega')$ for all $\omega' \in \Omega_V$. In the present section we consider only Markov chains that are induced by parallel kernels and corresponding stationary initial distributions. In that case, the interaction has the following representation (cf. [4]):

$$I(p, K) = \sum_{\nu \in V} (H(X_{\nu, n+1} | X_{\nu, n}) - H(X_{\nu, n+1} | X_n)). \quad (7)$$

Here, for each unit ν the term $H(X_{\nu, n+1} | X_{\nu, n}) - H(X_{\nu, n+1} | X_n)$ measures the reduction of the uncertainty about its next state by the additional knowledge of the current states of the other units. In the following we show representative simulations of small complex systems with strong interaction. The simulations implement the usual Markov dynamics on a set of N binary units together with a random search scheme to optimize the stochastic interaction of the Markov chains. That is, the interaction measure, I , is computed with respect to the induced stationary probability distribution of a parallel Markov kernel, and starting from initial random values the kernel is iteratively perturbed such that I increases. Details of the simulations will be described elsewhere [5]. Here, the main focus is on the almost deterministic nature of the optimized systems. Because the calculation of I is algorithmically complex, simulations are restricted to small N .

Figure 1 displays two optimized systems for $N = 3$ which are typical also for larger systems. The interaction, I , comes near to its theoretical upper bound

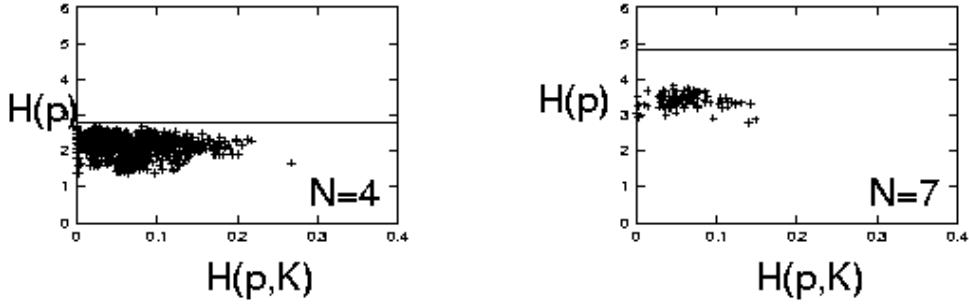


Fig. 2. Left: $H(p)$ and $H(p, K)$ for a series of simulations with $N = 4$ and $N = 7$. A theoretical upper bound for I and $H(p)$ is $N \ln 2 \approx 2.77$ and 4.85 , respectively.

of $N \ln 2 \approx 2.07$ in both examples, because the marginal kernel entropies are all near their maximum of $\ln 2 \approx .69$, and the full kernel entropies are very small (cf. (6)). This indicates that given the current state ω , the next state can be quite reliably predicted, but in contrast, not much information is gained from knowledge about the state of any individual unit alone. The Markov matrices in Fig. 1 (circle area represents transition probability) and derived state transition graphs (node-labels indicate states, ω , edge-labels transition probabilities) exhibit the determinism even more clearly: Most columns in the Markov matrices have only a single entry of probability 1 and, accordingly, most nodes in the transition graphs have just one outgoing edge.

The transition graphs reveal that the dynamics consists of an ergodic component comprising cycles of deterministic transitions nested by *branching states* (state 101 on the left, and state 001 on the right). Also visible are *transient states*, which, once left, are never occupied again (left: 100; right: 110) – thus, $p(\omega) = 0$ for transient states. Theorem 1 bounds the maximum number of outgoing transitions to $\leq N + 1$ for nodes with $p(\omega) > 0$ (the ergodic component). In a large number of simulations with N up to 8, we observed an increasing complexity in the transition graph structure, but never branching nodes with more than 2 outgoing transitions, although those should in principle be possible. In contrast, transient states can be deterministic (Fig. 1 left, node 100) as well as project to an arbitrary number of targets up to all possible ones (Fig. 1 right, node 110). Branching nodes and nested loops result in activity patterns that randomly switch between repetitive deterministic configuration sequences of various lengths (Fig. 1 lower frames).

Figure 2 shows distributions of $H(p)$ and $H(p, K)$ for $N = 4$ and 7. $H(p, K)$ is always small compared to the theoretical upper bound $N \ln 2$. So, the systems are almost deterministic. The ergodic component as in Fig. 1 consists of nested loops, but the graph structure gets increasingly complex and cycles long with system size (not shown). Also, the size of the transient component increases. Therefore, $H(p)$ falls below its maximal value of $N \ln 2$ for larger N .

Overall, Theorem 1 and the displayed simulations show that the maximization of stochastic interaction leads to Markov chains with highly deterministic global dynamics, but randomness in single unit activity. Future work will address stochastic interaction in more realistic neural network models.

References

- [1] Abeles, M. (1991) *Corticonics: Neural circuits of the cerebral cortex*. Cambridge: Cambridge University Press.
- [2] Amari, S.-I. (2001) Information Geometry on Hierarchy of Probability Distributions. *IEEE Transactions on Information Theory*, 47, 1701–1711.
- [3] Ay, N. (2001) An Information Geometric Approach to a Theory of Pragmatic Structuring. *The Annals of Probability*, in press.
- [4] Ay, N. (2002) Information Geometry on Complexity and Stochastic Interaction. Submitted.
- [5] Ay, N. and Wennekers, T. (2002) Dynamics of strongly interacting Markov chains. Submitted.
- [6] Eckhorn, R. (1999) Neural mechanisms of scene segmentation: Recordings from the visual cortex suggest basic circuits for linking field models. *IEEE Transactions on Neural Networks*, 10, 464–479.
- [7] Linsker, R. (1986) From Basic Network Principles to Neural Architecture. *Proceedings of the National Academy of Sciences*, USA 83, 7508–7512.
- [8] Rieke, F., Warland, D., Ruyter van Steveninck, R., & Bialek W. (1998) *Spikes: Exploring the Neural Code*. Cambridge: MIT Press.
- [9] Singer, W., & Gray, C.M. (1995) Visual feature integration and the temporal correlation hypotheses. *Annual Review of Neuroscience*, 18, 555–586.
- [10] Tononi, G., Sporns, O., & Edelman, G.M. (1994) A measure for brain complexity: Relating functional segregation and integration in the nervous system. *Proceedings of the National Academy of Sciences USA*, 91, 5033–5037.
- [11] Wennekers, T., & Palm G. (2000) Cell Assemblies, Associative Memory and Temporal Structure in Brain Signals. In: Miller, R. (ed.) *Time and the Brain. Conceptual Advances in Brain Research*, vol.2, pp. 251–273, Harwood Academic Publishers.