From finite-system entropy to entropy rate for a Hidden Markov Process

Or Zuk, Eytan Domany, Ido Kanter and Michael Aizenman

Abstract—A recent result presented the expansion for the entropy rate of a Hidden Markov Process (HMP) as a power series in the noise variable ϵ . The coefficients of the expansion around the noiseless ($\epsilon = 0$) limit were calculated up to 11th order, using a conjecture that relates the entropy rate of a HMP to the entropy of a process of finite length (which is calculated analytically). In this communication we generalize and prove the validity of the conjecture, and discuss the theoretical and practical consequences of our new theorem.

Index Terms-Hidden Markov Processes, Entropy rate

I. INTRODUCTION

ET $\{X_N\}$ be a finite state stationary Markov process • over the alphabet $\Sigma = \{1, .., s\}$, and let $\{Y_N\}$ be its noisy observation (on the same alphabet). The pair can be described by the Markov transition matrix $M = M_{s \times s}$ $\{m_{ij}\}\$ and the emission matrix $R = R_{s \times s}$, which yield the probabilities $P(X_{N+1} = j | X_N = i) = m_{ij}$ and $P(Y_N = i)$ $j|X_N = i) = r_{ij}$. We consider here the case where the signal to noise ratio (SNR) is small and M is strictly positive $(m_{ij} >$ 0) and thus has a unique stationary distribution. For the 'high - SNR' regime one may write $R = I + \epsilon T$, where $\epsilon > 0$ is some small number, I is the identity matrix, and the matrix $T = \{t_{ij}\}$ satisfies $t_{ii} < 0, t_{ij} \ge 0, \forall i \ne j \text{ and } \sum_{j=1}^{s} t_{ij} = 0.$ The process Y can be viewed as an observation of X through a noisy channel. It is an example of a Hidden Markov Process (HMP), and is determined by the parameters M, T and ϵ . More generally, HMPs have a rich and developed theory, and enormous applications in various fields (see [1], [2]).

An important property of Y is its entropy rate. The Shannon entropy rate of a stochastic process ([3]) measures the amount of 'uncertainty per-symbol'. More formally, for $i \leq j$ let $[X]_i^j$ denote the vector $(X_i, ..., X_j)$. The entropy rate is defined as:

$$\bar{H}(Y) = \lim_{N \to \infty} \frac{H([Y]_1^N)}{N} \tag{1}$$

Where $H(X) = -\sum_{X} P(X) \log P(X)$; We will sometimes omit the realization x of the variable X, so P(X) should

E.D. and O.Z. are in the Department of Physics of Complex Systems, Weizmann Inst. of Science.

I.K. is in the Department of Physics, Bar-Ilan Univ.

M.A. is in the Departments of Physics and Mathematics, Princeton Univ.

be understood as P(X = x). For a stationary process the limit in (1) exists and \overline{H} can also be computed via the conditional entropy ([4]) as: $\overline{H}(Y) = \lim_{N \to \infty} H(Y_N | [Y]_1^{N-1}).$ Here H(U|V) represents the conditional entropy, which for random variables U and V is the average uncertainty of the conditional distribution of U conditioned on V, that is $H(U|V) = \sum_{v} P(U = u) H(U|V = v)$. By the chain rule for entropy, it can also be viewed as a difference of entropies, H(U|V) = H(U, V) - H(V). This relation will be used below. There is at present no explicit expression for the entropy rate of a HMP ([1], [5]). Few recent works ([5], [6], [7]) have dealt with finding the asymptotic behavior of \overline{H} in several regimes, albeit giving rigorously only bounds or at most second ([7]) order behavior. Here we generalize and prove a relationship, that was posed in [7] as a conjecture, thereby turning the computation presented there, of H as a series expansion up to 11th order in ϵ , into a rigorous statement.

II. THEOREM STATEMENT AND PROOF

Our main result is the following:

Theorem 1: Let $H_N \equiv H_N(M, T, \epsilon) = H([Y]_1^N)$ be the entropy of a system of length N, and let $C_N = H_N - H_{N-1}$. Assume¹ there is some (complex) neighborhood $B_\rho(0) \subset \mathbb{C}$ of zero, in which the (one-variable) functions $\{C_N\}$ and \overline{H} are analytic in ϵ , with a Taylor expansion given by:

$$C_N(M,T,\epsilon) = \sum_{k=0}^{\infty} C_N^{(k)} \epsilon^k, \quad \bar{H}(M,T,\epsilon) = \sum_{k=0}^{\infty} C^{(k)} \epsilon^k$$
(2)

(The coefficients $C_N^{(k)}$ are functions of the parameters M and T. From now on we omit this dependence). Then:

$$N \ge \lceil \frac{k+3}{2} \rceil \Rightarrow C_N^{(k)} = C^{(k)} \tag{3}$$

 C_N is an upperbound ([4]) for \overline{H} . The behavior stated in Thm. 1 was discovered using symbolic computations, but proven only for $k \leq 2$, in the binary symmetric case ([7]). Although technically involved, our proof is based on two simple ideas.

¹It is easy to show that the functions C_N are differentiable to all orders in ϵ , at $\epsilon = 0$. The assumption which is not proven here is that they are in fact analytic with a radius of analyticity which is uniform in N, and are uniformly bounded within some common neighborhood of $\epsilon = 0$

First, we distinguish between the noise parameters at different sites. We consider a more general process $\{Z_N\}$, where Z_i 's emission matrix is $R_i = I + \epsilon_i T$. The process $\{Z_N\}$ is determined by M,T and $[\epsilon]_1^N$. We define the following functions:

$$F_N(M, T, [\epsilon]_1^N) = H([Z]_1^N) - H([Z]_1^{N-1})$$
(4)

Setting all the ϵ_i 's equal reduces us back to the Y process, so in particular $F_N(M, T, (\epsilon, ..., \epsilon)) = C_N(\epsilon)$.

Second, we observe that if a particular ϵ_i is set to zero, the observation Z_i equals the state X_i . Thus, conditioning back to the past is 'blocked'. This can be used to prove:

Lemma 1: Assume $\epsilon_j = 0$ for some 1 < j < N. Then:

$$F_N([\epsilon]_1^N) = F_{N-j+1}([\epsilon]_{j+1}^N)$$

Proof: F can be written as the sum:

$$F_N = -\sum_{[Z]_1^N} P([Z]_1^{N-1}) P(Z_N | [Z]_1^{N-1}) \log P(Z_N | [Z]_1^{N-1})$$
(5)

Here the dependence on $[\epsilon]_1^N$ and M, T is hidden in the probabilities P(..). Since $\epsilon_j = 0$, we must have $X_j = Z_j$, and therefore (since X is a Markov chain), conditioning further to the past is 'blocked', that is:

$$\epsilon_j = 0 \Rightarrow P(Z_N | [Z]_1^{N-1}) = P(Z_N | [Z]_j^{N-1})$$
 (6)

Substituting in eq. 5 gives:

$$F_{N} = -\sum_{[Z]_{1}^{N}} P([Z]_{1}^{N-1}) P(Z_{N}|[Z]_{j}^{N-1}) \log P(Z_{N}|[Z]_{j}^{N-1}) = -\sum_{[Z]_{j}^{N}} P([Z]_{j}^{N}) \log P(Z_{N}|[Z]_{j}^{N-1}) = F_{N-j+1}$$
(7)

Let $\vec{k} = [k]_1^N$ be a vector with $k_i \in \{\mathbb{N} \cup 0\}$. Define its 'weight' as $\omega(\vec{k}) = \sum_{i=1}^N k_i$. Define also:

$$F_N^{\vec{k}} \equiv \left. \frac{\partial^{\omega(\vec{k})} F_N}{\partial \epsilon_1^{k_1}, .., \partial \epsilon_N^{k_N}} \right|_{\vec{\epsilon} = 0} \tag{8}$$

As we now show, adding zeros to \vec{k} leaves $F_N^{\vec{k}}$ unchanged :

Lemma 2: Let $\vec{k} = [k]_1^N$ with $k_1 \leq 1$. Denote $\vec{k}^{(r)}$ the concatenation: $\vec{k}^{(r)} = (0, ..., 0, k_1, ..., k_N)$. Then:

$$F_{N}^{\vec{k}} = F_{r+N}^{\vec{k}^{(r)}} \quad , \forall r \in \mathbb{N}$$

Proof: Assume first $k_1 = 0$. Using lemma 1, we get:

$$F_{N+r}^{\vec{k}^{(r)}}([\epsilon]_{1}^{N+r}) = \frac{\partial^{\omega(\vec{k}^{(r)})}F_{r+N}([\epsilon]_{1}^{N+r})}{\partial\epsilon_{r+2}^{k_{2}}, ..., \partial\epsilon_{r+N}^{k_{N}}}\bigg|_{\vec{\epsilon}=0} = \frac{\partial^{\omega(\vec{k})}F_{N}([\epsilon]_{r+1}^{N+r})}{\partial\epsilon_{r+2}^{k_{2}}, ..., \partial\epsilon_{r+N}^{k_{N}}}\bigg|_{\vec{\epsilon}=0} = F_{N}^{\vec{k}}([\epsilon]_{r+1}^{r+N})$$
(9)

The case $k_1 = 1$ is reduced back to the case $k_1 = 0$ by taking the derivative. We denote $[Z]_1^{N(j \to r)}$ the vector which is equal to $[Z]_1^N$ in all coordinates except on coordinate j, where $Z_j = r$. Using eq. 9, we get:

$$\begin{split} F_{N+1}^{\vec{k}^{(1)}}([\epsilon]_{1}^{N+1}) &= \frac{\partial^{\omega(k)-1}}{\partial \epsilon_{3}^{k_{2}} \dots \partial \epsilon_{N+1}^{k_{N}}} \left[\frac{\partial F_{N+1}}{\partial \epsilon_{2}} \Big|_{\epsilon_{2}=0} \right] \Big|_{\vec{\epsilon}=0} = \\ & \frac{\partial^{\omega(\vec{k})-1}}{\partial \epsilon_{3}^{k_{2}} \dots \partial \epsilon_{N+1}^{k_{N}}} \left\{ -\sum_{r=1}^{s} t_{X_{i}r} \sum_{[Z]_{1}^{N+1}} \right. \\ & \left[P([Z]_{1}^{N+1^{(2\to r)}}) \log P(Z_{N+1}|[Z]_{1}^{N}) - \right. \\ & \left. P(Z_{N+1}|[Z]_{1}^{N}) P([Z]_{1}^{N^{(2\to r)}}) \right] \Big|_{\epsilon_{2}=0} \right\} \Big|_{[\epsilon]_{1}^{N+1}=0} = \\ & \left. \frac{\partial^{\omega(\vec{k})-1}}{\partial \epsilon_{2}^{k_{2}} \dots \partial \epsilon_{N}^{k_{N}}} \left\{ -\sum_{r=1}^{s} t_{X_{i}r} \sum_{[Z]_{1}^{N}} \right. \\ & \left. \left[P([Z]_{1}^{N^{(1\to r)}}) \log P(Z_{N}|[Z]_{1}^{N-1}) - \right. \\ \\ & \left. P(Z_{N}|[Z]_{1}^{N-1}) P([Z]_{1}^{N^{(1\to r)}}) \right] \Big|_{\epsilon_{1}=0} \right\} \Big|_{[\epsilon]_{1}^{N}=0} = F_{N}^{\vec{k}}([\epsilon]_{1}^{N}) \end{split}$$

 $C_N^{(k)}$ is obtained by summing $F_N^{\vec{k}}$ on all \vec{k} 's with weight k:

$$C_N^{(k)} = \sum_{\vec{k}, \omega(\vec{k}) = k} F_N^{\vec{k}}$$
(11)

The next lemma shows that one does not need to sum on all such \vec{k} 's, as many of them give zero contribution:

Lemma 3: Let $\vec{k} = [k]_1^N$. If $\exists i < j < N$, with $k_i \ge 1, k_j \le 1$, then $F_N^{\vec{k}} = 0$.

Proof: Assume first $k_j = 0$. Using lemma 1 we get

$$F_{N}^{\vec{k}} \equiv \left. \frac{\partial^{\omega(\vec{k})} F_{N}(\vec{\epsilon})}{\partial \epsilon_{1}^{k_{1}}, ..., \partial \epsilon_{N}^{k_{N}}} \right|_{\vec{\epsilon}=0} = \left. \frac{\partial^{\omega(\vec{k})} F_{N-j+1}([\epsilon]_{j}^{N})}{\partial \epsilon_{1}^{k_{1}}, ..., \partial \epsilon_{N}^{k_{N}}} \right|_{\vec{\epsilon}=0} = \left. \frac{\partial^{\omega(\vec{k})-1}}{\partial \epsilon_{1}^{k_{1}}, ..., \partial \epsilon_{N}^{k_{i}-1}} \right|_{\vec{\epsilon}=0} = 0 \quad (12)$$

Assume now $k_j = 1$. Write the probability of Z:

$$P([Z]_{1}^{N}) = \sum_{[X]_{1}^{N}} P([X]_{1}^{N}) P([Z]_{1}^{N} | [X]_{1}^{N}) = \sum_{[X]_{1}^{N}} P([X]_{1}^{N}) \prod_{i=1}^{N} (\delta_{X_{i}Z_{i}} + \epsilon_{i}t_{X_{i}Z_{i}})$$
(13)

where δ is Kronecker's delta. Differentiate with respect to ϵ_j :

$$\left. \frac{\partial P([Z]_1^N)}{\partial \epsilon_j} \right|_{\epsilon_j = 0} =$$

(10)

$$\sum_{[X]_1^N} \left[P([X]_1^N) t_{X_j Z_j} \prod_{i \neq j} (\delta_{X_i Z_i} + \epsilon_i t_{X_i Z_i}) \right] \bigg|_{\epsilon_j = 0} = \left\{ \left. \left\{ \sum_{r=1}^s t_{X_i r} P([Z]_1^N^{(j \to r)}) \right\} \right|_{\epsilon_j = 0}$$
(14)

Using Bayes' rule $P(Z_N | [Z]_1^{N-1}) = \frac{P([Z]_1^N)}{P([Z]_1^{N-1})}$, we get:

$$\frac{\partial P(Z_N | [Z]_1^{N-1})}{\partial \epsilon_j} \bigg|_{\epsilon_j = 0} =$$

$$\frac{1}{P([Z]_1^{N-1})} \sum_{r=1}^s t_{X_i r} \left[P([Z]_1^{N^{(j \to r)}}) - P([Z]_1^{N-1}) P([Z]_1^{N-1^{(j \to r)}}) \right] \bigg|_{\epsilon_j = 0}$$
(15)

This gives:

$$\frac{\partial [P([Z]_1^N) \log P(Z_N | [Z]_1^{N-1})]}{\partial \epsilon_j} \bigg|_{\epsilon_j = 0} = \sum_{r=1}^s t_{X_i r} \left\{ P([Z]_1^{N(j \to r)}) \log P(Z_N | [Z]_1^{N-1}) + \right\}$$

$$P([Z]_{1}^{N^{(j\to r)}}) - P(Z_{N}|[Z]_{1}^{N-1})P([Z]_{1}^{N-1^{(j\to r)}}) \bigg\} \bigg|_{\epsilon_{j}=0}$$
(16)

And therefore:

$$\frac{\partial F_N}{\partial \epsilon_j} \bigg|_{\epsilon_j = 0} = -\sum_{r=1}^s t_{X_i r} \left\{ \sum_{[Z]_1^N} \left[P([Z]_1^{N(j \to r)}) \log P(Z_N | [Z]_1^{N-1}) - P(Z_N | [Z]_1^{N-1}) P([Z]_1^{N-1})^{(j \to r)}) \right] \right\} \bigg|_{\epsilon_j = 0} = \left\{ -\sum_{r=1}^s t_{X_i r} \sum_{[Z]_j^N} \left[P([Z]_j^{N(1 \to r)}) \log P(Z_N | [Z]_j^{N-1}) - P(Z_N | [Z]_j^{N-1}) P([Z]_j^{N-1})^{(1 \to r)}) \right] \right\} \bigg|_{\epsilon_1 = 0}$$

$$\left\{ -\left| \sum_{r=1}^s T_{X_i r} \sum_{[Z]_j^N} \left[P([Z]_j^{N-1}) P([Z]_j^{N-1}) - P(Z_N | [Z]_j^{N-1}) \right] \right\} \bigg|_{\epsilon_1 = 0}$$

$$\left\{ -\left| \sum_{r=1}^s T_{X_i r} \sum_{[Z]_j^N} \left[P([Z]_j^{N-1}) P([Z]_j^{N-1}) - P(Z_N | [Z]_j^{N-1}) \right] \right\} \right|_{\epsilon_1 = 0}$$

$$\left\{ -\left| \sum_{r=1}^s T_{X_i r} \sum_{[Z]_j^N} \left[P([Z]_j^{N-1}) P([Z]_j^{N-1}) - P(Z_N | [Z]_j^{N-1}) \right] \right\} \right|_{\epsilon_1 = 0}$$

$$\left\{ -\left| \sum_{r=1}^s T_{X_i r} \sum_{[Z]_j^N} \left[P([Z]_j^{N-1}) P([Z]_j^{N-1}) - P(Z_N | [Z]_j^{N-1}) \right] \right\} \right\} \right\} \right\}$$

The latter equality comes from using eq. 6, which 'blocks' the dependence backwards. Eq. 17 shows that ϵ_i does not appear in $\frac{\partial F_N}{\partial \epsilon_j}\Big|_{\epsilon_j=0}$ for i < j, therefore $\frac{\partial^{k_i+1}F_N}{\partial \epsilon_i^{k_i}\partial \epsilon_j}\Big|_{\epsilon_j=0} = 0$ and $F_N^{\vec{k}} = 0$.

We are now ready to prove our main theorem:

Proof:

Let $\vec{k} = [k]_1^N$ with $\omega(\vec{k}) = k$. Define its 'length' as $l(\vec{k}) = N + 1 - \min_{k_i > 1}\{i\}$. It easily follows from lemma 3 that $F_N^{\vec{k}} \neq 0 \Rightarrow l(\vec{k}) \leq \lceil \frac{k+3}{2} \rceil - 1$. Thus, according to lemma 2:

$$F_{N}^{\vec{k}} = F_{\lceil \frac{k+3}{2} \rceil}^{(k_{N-\lceil \frac{k+3}{2} \rceil+1}, \dots, k_{N})}$$
(18)

for all \vec{k} 's in the sum. Summing on all $F_N^{\vec{k}}$ with the same 'weight' gives $C_N^{(k)} = C_{\lceil \frac{k+3}{2} \rceil}^{(k)}$, $\forall N > \lceil \frac{k+3}{2} \rceil$. But from the analyticity of C_N and H near $\epsilon = 0$ it follows that $\lim_{N\to\infty} C_N^{(k)} = C^{(k)}$, therefore $C_N^{(k)} = C^{(k)}$, $\forall N \ge \lceil \frac{k+3}{2} \rceil$.

III. CONCLUSION

Our main theorem sheds light on the connection between finite and infinite chains, and gives a practical and straightforward way to compute the entropy rate as a series expansion in ϵ up to an arbitrary power. The surprising 'settling' of the expansion coefficients $C_N^{(k)} = C^{(k)}$ for $N \ge \lceil \frac{k+3}{2} \rceil$, hold for the entropy. For other functions involving only conditional probabilities (e.g. relative entropy between two *HMPs*) a weaker result holds: the coefficients 'settle' for $N \ge k$. One can expand the entropy rate in several parameter regimes. As it turns out, exactly the same 'settling' as was proven in Thm. 1 happens in the 'almost memoryless' regime, where M is close to a matrix which makes the X_i 's i.i.d. This and other regimes, as well as the analytic behavior of the *HMP* ([8]), will be discussed elsewhere.

ACKNOWLEDGMENT

M.A. is grateful for the hospitality shown him at the Weizmann Institute, where his work was supported by the Einstein Center for Theoretical Physics and the Minerva Center for Nonlinear Physics. The work of I.K. at the Weizmann Institute was supported by the Einstein Center for Theoretical Physics. E.D. and O.Z. were partially supported by the Minerva Foundation and by the European Community's Human Potential Programme under contract HPRN-CT-2002-00319, STIPCO.

REFERENCES

- Y. Ephraim and N. Merhav, *Hidden Markov processes*, IEEE Trans. Inform. Theory, 48(6), pp. 1518-1569, 2002.
- [2] L. R. Rabiner, A tutorial on hidden Markov models and selected applications in speech recognition, Proc. IEEE, 77(2), pp. 257-286, 1989.
- [3] C. E. Shannon, A mathematical theory of communication, Bell System Technical Journal, 27, pp. 379-423 and 623-656, 1948.
- [4] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, Wiley, New York, 1991.
- [5] P. Jacquet, G. Seroussi and W. Szpankowski, On the Entropy of a Hidden Markov Process, DCC 2004, pp. 362-371.

- [6] E. Ordentlich and T. Weissman, New Bounds on the Entropy Rate of Hidden Markov Processes, San Antonio Information Theory Workshop, 2004.
- [7] O. Zuk, I. Kanter and E. Domany, *Asymptotics of the Entropy Rate for a Hidden Markov Process*, DCC 2005, pp. 173-182.
- [8] G. Han and B. Marcus Analyticity of Entropy Rate in Families of Hidden Markov Chains, to appear in ISIT 2005