

Complex Systems Methods — 3. Statistical complexity of temporal sequences

Eckehard Olbrich

MPI MiS Leipzig

Potsdam WS 2007/08

- 1 Summary Kolmogorov sufficient statistic
- 2 Intuitive notions of complexity
- 3 Statistical complexity
 - Entropy convergence and excess entropy
 - Predictive information
- 4 Entropy estimation
 - Entropy of a discrete random variable — finite sample corrections

Summary Kolmogorov sufficient statistic

- Kolmogorov complexity

$$K_{\mathcal{U}}(x|I(x)) = \min_{p:\mathcal{U}(p,I(x))=x} I(p)$$

- Kolmogorov structure function, x^n denotes a string of length n

$$K_k(x^n|n) = \min_{\substack{p : I(p) \leq k \\ \mathcal{U}(p, n) = S \\ x^n \in S \subseteq \{0, 1\}^n}} \log |S|$$

- Kolmogorov sufficient statistic: least k such that

$$K_k(x^n|n) + k \leq K(x^n|n) + c .$$

- Regularities in x are described by describing the **set** S . Given S the string x is random,
- Algorithmic complexity \equiv randomness. But intuitively, complexity measure should quantify structure, not randomness. Thus it is related to an ensemble and not a single object.

all continuous	Partial Differential Equations (PDE's)
discretized space	coupled ordinary differential equations (ODE's)
discretized time	coupled map lattice (CML)
discretized state	cellular automata (CA)

- all dynamical systems either deterministic or stochastic
- Digital computer: Finite state automaton - finite number of discrete states

Intuitive notions of complexity — Stephen Wolframs classification of cellular automata

Elementary cellular automata: binary states $\{0, 1\}$. Next neighbour interaction. Rules can be coded by numbers $0, \dots, 255$:

111	110	101	100	011	010	001	000	
0	0	0	1	1	1	1	0	e.g. Rule 30.

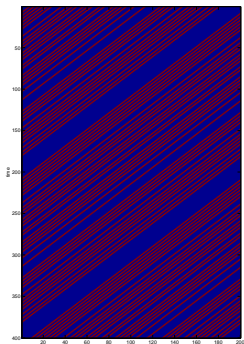
Class 1: Evolution leads to a homogeneous state.

Class 2: Evolution leads to a set of separated simple stable or periodic structures.

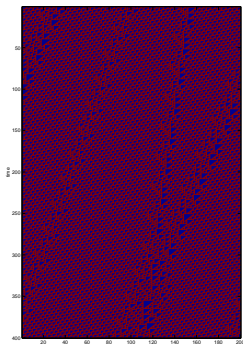
Class 3: Evolution leads to a chaotic pattern.

Class 4: Evolution leads to complex localized structures, sometimes long-lived.

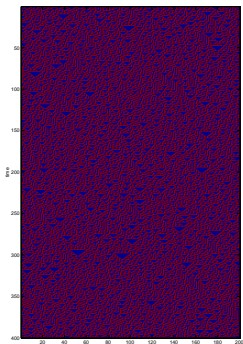
Intuitive notions of complexity — Stephen Wolfram's classification of cellular automata



class 2



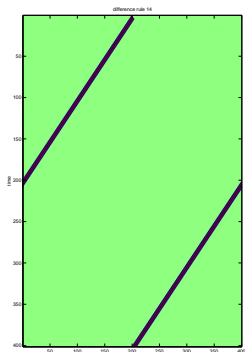
class 4



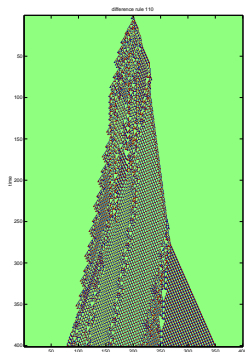
class 3

Intuitive notions of complexity — Stephen Wolfram's classification of cellular automata

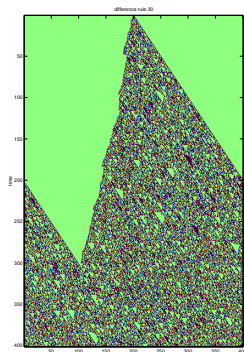
Propagation of a single perturbation



class 2
ordered



class 4
"complex"



class 3
chaotic

Statistical complexity — excess entropy

- $K(x|I(x))$ — minimal length of a program, which produces exactly the string x . Most complex strings are algorithmically random.
- Kolmogorov sufficient statistic: program that describes all regularities in x and computes the set S , which contains all strings with the same regularities as in x , but are otherwise random. Algorithmic complexity can be divided in two parts - regularities and randomness.
- Same idea for a time series (infinite sequence) provided with a stationary distribution:
 - Randomness per symbol is given by the entropy rate

$$h_\infty = \lim_{n \rightarrow \infty} h_n \quad h_n = H(X_0 | X_{-1}, \dots, X_{-n})$$

- Regularities quantified by the excess entropy (Crutchfield) or effective measure complexity (Grassberger)

$$E = \lim_{n \rightarrow \infty} E_n \quad \text{with} \quad E_n = (H_n - n \cdot h_n) \quad \text{and} \quad H_n = H(X_n, \dots, X_1)$$

Entropy convergence and excess entropy

The idea which led originally to this complexity measure was to use the converge rate of the conditional entropy to quantify the complexity of a time series:

fast convergence \rightarrow short memory \rightarrow low complexity

slow convergence \rightarrow long memory \rightarrow high complexity.

Conditional entropies

$$h_n = H(X_0|X_{-1}, \dots, X_{-n}) = H(X_0, X_{-1}, \dots, X_{-n}) - H(X_{-1}, \dots, X_{-n})$$

decrease monotonically. Decay is quantified by the conditional mutual information

$$\delta h_n := h_{n-1} - h_n = MI(X_0 : X_{-n}|X_{-1}, \dots, X_{-n}) .$$

Entropy convergence and excess entropy

Excess entropy:

$$\begin{aligned} E_N &= H(X_1, \dots, X_N) - N h_N \\ &= \sum_{n=0}^{N-1} (h_n - h_N) \quad \text{using } h_0 = H(X_1) \\ &= \sum_{n=0}^{N-1} \sum_{k=n+1}^N \delta h_k \quad \text{using } h_{n-1} = \delta h_n + h_n \\ &= \sum_{k=1}^N k \delta h_k \end{aligned}$$

The slower the convergence of the entropy the larger the excess entropy. If the sequence is Markov' of order m : $\delta h_n = 0$ for $n > m$. Thus $E = E_m$.

Bialek et al. (2000) proposed to quantify the complexity of a time series by the amount of information which the past tells us about the future

$$I_{pred} = MI(\text{past} : \text{future}) = \lim_{n_p, n_f \rightarrow \infty} MI(X_{-n_p}, \dots, X_0 : X_1, \dots, X_{n_f})$$

and called this *predictive information*.

The predictive information is equal to the excess entropy if the corresponding limits exist:

$$I_{pred} = E$$

$$\begin{aligned}I_{pred} &= \lim_{n_p, n_f \rightarrow \infty} MI(X_{-n_p}, \dots, X_{-1}, X_0 : X_1, X_2, \dots, X_{n_f}) \\&= \lim_{n_p, n_f \rightarrow \infty} \{H(X_1, \dots, X_{n_f}) + H(X_{-n_p}, \dots, X_0) \\&\quad - H(X_{-n_p}, \dots, X_0, \dots, X_{n_p})\} \\&= \lim_{n_p, n_f \rightarrow \infty} \{E_{n_f} + n_f h_{n_f} + E_{n_p+1} + (n_p + 1)h_{n_p+1} \\&\quad - E_{n_p+n_f+1} - (n_f + n_p + 1)h_{n_f+n_p+1}\} \\&= E\end{aligned}$$

The excess entropy measures the amount of information which is available from the past for predicting the time series.

Causal equivalence: Two past sequences $x_{-\infty}^0 = x_0, x_{-1}, \dots$ and $x_{-\infty}'^0 = x'_0, x'_{-1}, \dots$ are causal equivalent, if they have the same future, i.e. $p(x_1^\infty | x_{-\infty}^0) = p(x_1^\infty | x_{-\infty}'^0)$. The equivalence classes are called *causal states*. This notion was introduced by James P. Crutchfield et al. They called the transition graph between the causal states an *ϵ -machine*. In general it seems to be a subclass of hidden Markov models. Crutchfield and Young (1989): Statistical complexity C_μ as the entropy of the stationary distribution over the states. In general there is

$$C_\mu \geq E .$$

Grassberger (1986) called it *set complexity* in the context of regular languages.

In general: The *excess entropy* provides a **lower bound** for the amount of information necessary for an **optimal prediction**.

Excess entropy — some examples

- Sequence with period n : $h_m = 0$ for $m \geq n$, in particular $h_\infty = 0$. The entropy $H_{m \geq n} = \log n$ remains constant for $m \geq n$, thus $E = \log n$. All periodic sequences of the same length have the same complexity according to this measure. (Alternative: transient information introduced by Crutchfield and Feldman).
- Markov chain: $h_\infty = h_1$. $E = \delta h_1 = h_0 - h_1 = MI(X_{n+1} : X_n)$.
- Chaotic maps: Usually exponential decay of h_n — finite E .
- Feigenbaum point: $h_\infty = 0$, but $h_n \propto 1/n$, thus E diverges.

Complexity for the period doubling route to chaos

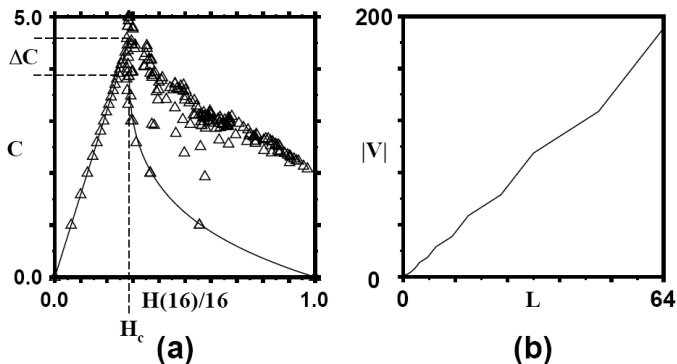


Figure 6 (a) Statistical complexity C_ν versus specific entropy $H(L)/L$ for the period-doubling route to chaos. Triangles denote estimated $(C_\nu, H(L)/L)$ at 193 values of the logistic map nonlinearity parameter. ϵ -machines were reconstructed using a subsequence length of $L = 16$. The heavy solid lines overlaying some of this empirical data are the analytical curves derived for C_0 versus $H_0(L)/L$. (After [24].) (b) At one of the critical parameter values of the period-doubling cascade in the logistic map the number $\|V\|$ of inferred states grows without bound. Here $r = r_c \approx 3.5699456718695445 \dots$ and the sequence length ranges up to $L = 64$ where $\|V\| = 196$ states are found. It can be shown, and can be inferred from the figure, that the per symbol density of states $\|V(L)\|/L$ does not have a limiting value as $L \rightarrow \infty$. (After [56].)

From Crutchfield/Young 1990. Computation at the onset of chaos.

Excess entropy and attractor dimension

Deterministic system with continuous state observables: The entropy $H_m(\epsilon) = H(X_1, \dots, X_m; \epsilon)$ scales with respect to the resolution ϵ :

$$m < D$$

$$H_m(\epsilon) = H_m^c - m \log \epsilon + \mathcal{O}(\epsilon)$$

$$m > D$$

$$H_m(\epsilon) = \text{const} - D \log \epsilon + \mathcal{O}(\epsilon)$$

Because $h_\infty = h_{KS}$ does not depend on ϵ we have

$$E \propto -D \log \epsilon .$$

The better one knows the initial conditions of the system the better the system can be predicted.

- Peter Grassberger, Toward a quantitative theory of self-generated complexity, *International Journal of Theoretical Physics*, 25 (1986), 907-938.
- James P. Crutchfield and David P. Feldman, Regularities unseen, randomness observed: Levels of entropy convergence *Chaos*, 13 (2003), 25-54.
- William Bialek, Ilya Nemenman and Naftali Tishby, Predictability, complexity, and learning. *Neural Computation*, 13 (2001), 2409-2463.
- Remo Badii and Antonio Politi, *Complexity — Hierarchical structures and scaling in physics*, Cambridge University Press, 1997.

Entropy of a discrete random variable — finite sample corrections

Reference: P. Grassberger, Entropy Estimates from Insufficient Samplings, arXiv:physics/0307138v1

- N data points randomly and *independently* distributed on M boxes. The number n_i of points in each box is a random variable with an expectation value $z_i := E[n_i] = p_i N$. Their distribution is binomial

$$P(n_i; p_i, N) = \binom{N}{n_i} p_i^{n_i} (1 - p_i)^{N - n_i}$$

- For $p_i \ll 1 \quad \forall i$ the n_i can be assumed to be Poisson distributed

$$P(n_i; z_i) = \frac{z_i^{n_i}}{n_i!} e^{-z_i}$$

Entropy of a discrete random variable — finite sample corrections

- Entropy

$$H = - \sum_{i=1}^M p_i \log p_i = \ln N - \frac{1}{N} \sum_{i=1}^M z_i \ln z_i$$

- Naive estimator

$$\hat{H}_{naive} = \ln N - \frac{1}{N} \sum_{i=1}^M n_i \ln n_i$$

- In general the estimator is biased, i.e.

$$\Delta H := E[\hat{H}] - H \neq 0. \quad \Delta H_{naive} < 0.$$

Entropy of a discrete random variable — finite sample corrections

- In the limit of large N and M each contribution $z_i \ln z_i$ will be statistically independent and can be estimated as function of n_i :

$$z_i \ln z_i \approx z_i \hat{\ln} z_i = n_i \phi(n_i) \quad E[z_i \hat{\ln} z_i] = \sum_{n_i=1}^{\infty} n_i \phi(n_i) P(n_i; z_i) .$$

Implizit assumption: $n_i = 0$ gives no information about p_i .

- Estimator

$$\hat{H}_\phi = \ln N - \frac{1}{N} \sum_{i=1}^M n_i \phi(n_i)$$

- Grassbergers result:

$$E[n\psi(n)] = z \ln z + zE_1(z)$$

with the digamma function

$$\psi(x) = \frac{d \ln \Gamma(x)}{dx} \quad \text{and} \quad E_1 = \Gamma(0, x) = \int_1^\infty \frac{e^{-xt}}{t} dt .$$

- For large z $zE_1(z) \approx e^{-z}$, thus neglecting this term gives the estimator

$$\hat{H}_\psi = \ln N - \frac{1}{N} \sum_{i=1}^M n_i \psi(n_i) .$$

- Approximations for the digamma function leads to known estimators:
 - ① $\psi(x) \approx \ln x$: naive estimator
 - ② $\psi(x) \approx \ln x - 1/(2x)$: Miller correction
 - ③ $\psi(x) < \ln x - 1/(2x) < \ln x$ leads to $\hat{H}_\psi > \hat{H}_{Miller} > \hat{H}_{naive}$
- Grassbergers best estimator:

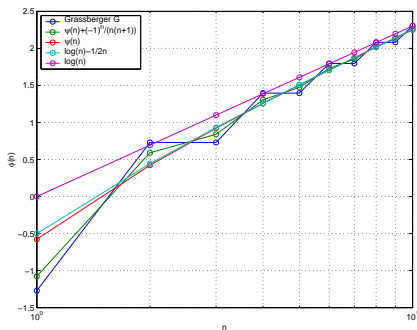
$$\hat{H}_G = \ln N - \frac{1}{N} \sum_{i=1}^M n_i G_{n_i}$$

with

$$G_n = \psi(n) + (-1)^n \sum_{k=1}^{\infty} \frac{1}{(n+2k)(n+2k+1)}$$

or $G_1 = -\gamma - \ln 2$ $G_2 = 2 - \gamma - \ln 2$ $G_{2n+1} = G_{2n}$ and
 $G_{2n+2} = G_{2n} + \frac{2}{2n+1}$ for $n \geq 1$.

Finite sample corrections — Comparison between different corrections

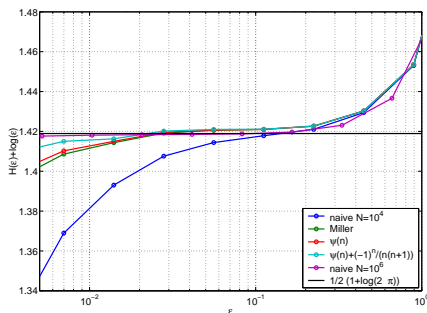


Note that $\psi(x) + \frac{(-1)^n}{n(n+1)}$ corresponds to an approximation of

$$G_n = \psi(n) + (-1)^n \sum_{k=1}^{\infty} \frac{1}{(n+2k)(n+2k+1)}$$

with considering only the $k = 0$ term in the sum.

Test — Entropy of a Gaussian distribution



Differential entropy of a Gaussian distribution

$$H_{Gauss}^C = \frac{1}{2}(1 + \log(2\pi\sigma^2)) = H(\epsilon) + \log(\epsilon) + \mathcal{O}(\epsilon)$$