

Entropy

Jonathan L.F. King
University of Florida, Gainesville FL 32611-2082, USA
squash@math.ufl.edu
Webpage <http://www.math.ufl.edu/~squash/>
4 May, 2008 (at 0-900:20)

Prolegomenon. The word ‘*entropy*’ was coined by Rudolf Julius Emanuel Clausius in 1867, in [2], referring to the thermodynamic notion in physics.

Our focus here, however, will be the notion in measurable-dynamics and topological-dynamics. (Entropy in differentiable-dynamics,^{♥1} would require an article by itself.) Shannon’s 1948 paper [3] on Information Theory, then Kolmogorov’s [4] and Sinai’s [5] generalization to dynamical systems, will be our starting point. I will stay in the one-dimensional case, where the acting-group is \mathbb{Z} .

§Bird’s-eye view

Prolegomenon	1	Cautions on determinism’s relation to zero-entropy	15
Glossary	1	Pinsker Field	15
Entropy example: <i>How many questions?</i>	3	Ornstein theory	16
Distribution entropy	3	The Pinsker-Field and K-automorphisms	16
The η function	4	Topological Entropy	17
Binomial coefficients	5	Using a metric	18
A gander at Shannon’s Noisy Channel theorem	6	Metric preliminaries	18
Noisy Channel	6	You take the High Road and I’ll take the Low Road	19
The information function	7	E.g: $\mathcal{E}_{\text{top}}(\text{Isometry}) = 0$	21
Entropy of a Process	9	E.g: <i>Topological Markov Shifts</i>	21
Bernoulli processes	10	The variational principle	24
Entropy of a Transformation	10	Topology on \mathfrak{M}	24
Entropy is continuous	11	Three recent results	27
Entropy is not continuous	11	Ornstein-Weiss: Finitely-observable in- variant	27
E.g: Meshalkin’s map	12	Entropy of actions of free groups	28
E.g: <i>Markov Shifts</i>	13	Conclusion	29
Determinism and Zero-entropy	14		
E.g: <i>Rotations are deterministic</i>	14		
E.g: <i>Rank-1 has zero-entropy</i>	15		

Glossary. Some of the following definitions refer to the “Notation” paragraph immediately below. Use *mpt* for ‘*measure-preserving transformation*’.

A **measure space** (X, \mathcal{X}, μ) is a set X , a **field** (that is, a σ -algebra) \mathcal{X} of subsets of X , and a countably-additive measure $\mu: \mathcal{X} \rightarrow [0, \infty]$. (We often just write (X, μ) , with the field implicit.) For a collection $\mathcal{C} \subset \mathcal{X}$, use $\text{Fld}(\mathcal{C})$ for the smallest field $\supset \mathcal{C}$. The number $\mu(B)$ is the “ μ -mass of B ”.

A **measure-preserving map** $\psi: (X, \mathcal{X}, \mu) \rightarrow (Y, \mathcal{Y}, \nu)$ is a map $\psi: X \rightarrow Y$ such that the inverse image of each $B \in \mathcal{Y}$ is in \mathcal{X} , and $\mu(\psi^{-1}(B)) = \nu(B)$. A **(measure-preserving) transformation** is a measure-preserving map $T: (X, \mathcal{X}, \mu) \rightarrow (X, \mathcal{X}, \mu)$. Condense this notation to $(T: X, \mathcal{X}, \mu)$ or $(T: X, \mu)$.

A **probability space** is a measure space (X, μ) with $\mu(X) = 1$; this μ is a **probability measure**. All our maps/transformations in this article are on probability spaces. A **factor map**

$$\psi: (T: X, \mathcal{X}, \mu) \rightarrow (S: Y, \mathcal{Y}, \nu)$$

^{♥1}For instance, see [24], [26], [18] and [15].

is a measure-preserving map $\psi: X \rightarrow Y$ which intertwines the transformations, $\psi \circ T = S \circ \psi$. And ψ is an **isomorphism** if –after deleting a nullset in each space– this ψ is a bijection and ψ^{-1} is also a factor map.

A measure-theoretic statement holds **almost everywhere**, abbreviated **a.e.**, if it holds off of a nullset, mass-zero set.^{♥2} For example, $B \supset_{\text{a.e.}} A$ means that $\mu(B \setminus A)$ is zero. The *a.e.* will usually be implicit.

A **probability vector** $\vec{v} = (v_1, v_2, \dots)$ is a list of non-negative reals whose sum is 1. We generally assume that probability vectors and partitions (see below) have *finitely* many components. We write “countable probability vector/partition”, when finitely or denumerably many components are considered.

A **partition** $P = (A_1, A_2, \dots)$ splits X into pairwise disjoint subsets $A_i \in \mathcal{X}$ so that the disjoint union $\bigsqcup_i A_i$ is all of X . Each A_i is an **atom** of P . Use $|P|$ or $\#P$ for the number of atoms. When P partitions a probability space, then it yields a probability vector \vec{v} , where $v_j := \mu(A_j)$. Lastly, use $P\langle x \rangle$ to denote the P -atom that owns x .

Fonts. We use the font \mathcal{H} , \mathcal{E} , \mathcal{I} for *distribution-entropy*, *entropy* and the *information function*. In contrast, the script font $\mathcal{ABC} \dots$ will be used for collections of sets; usually subfields of \mathcal{X} . Use $\mathbb{E}(\cdot)$ for the (conditional) expectation operator.

Notation. \mathbb{Z} = integers. \mathbb{Z}_+ = positive integers, and \mathbb{N} = natural numbers^{♥3} = $\{0, 1, 2, \dots\}$. Use $\lceil \cdot \rceil$ and $\lfloor \cdot \rfloor$ for the **ceiling** and **floor** functions; $\lfloor \cdot \rfloor$ is also called the “greatest-integer function”. For an interval $J := [a, b) \subset [-\infty, +\infty]$, let $[a .. b)$ denote the **interval of integers** $J \cap \mathbb{Z}$ (with a similar convention for closed and open intervals). E.g., $(e .. \pi] = (e .. \pi) \cap \mathbb{Z} = \{3\}$.

For subsets A and B of the same space, Ω , use $A \subset B$ for inclusion and $A \subsetneq B$ for *proper* inclusion. The difference set $B \setminus A$ is $\{\omega \in B \mid \omega \notin A\}$. Employ A^c for the complement $\Omega \setminus A$. Since we work in a probability space, if we let $x := \mu(A)$, then a convenient convention is to have

$$x^c := 1 - x,$$

since then $\mu(A^c)$ equals x^c .

Use $A \triangle B$ for the **symmetric difference** $[A \setminus B] \cup [B \setminus A]$. For a collection $\mathcal{C} = \{E_j\}_j$ of sets in Ω , let the **disjoint union** $\bigsqcup_j E_j$ or $\bigsqcup(\mathcal{C})$ represent the union $\bigcup_j E_j$ and also assert that the sets are pairwise disjoint.

Use “ $\forall_{\text{large } n}$ ” to mean: “ $\exists n_0$ such that $\forall n > n_0$ ”. To refer to lefthand side of an equation 17, use LhS(17); do analogously for RhS(17), the righthand side.

^{♥2}Eugene Gutkin once remarked to me that the problem with Measure Theory is... that you have to say “almost everywhere”, almost everywhere.

^{♥3}Some well-meaning folk use \mathbb{N} for \mathbb{Z}_+ , saying ‘*Nothing could be more natural than the positive integers*’. And this is why $0 \in \mathbb{N}$.

Entropy example: How many questions? Imagine a dartboard, FIG. 1, split in five regions, A, \dots, E , with known probabilities. Blindfolded, you throw a dart at the board. What is the expected number, V , of Yes/No questions needed to ascertain the region in which the dart landed?

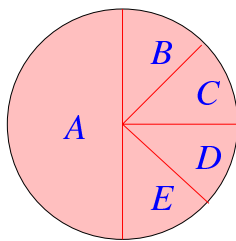


FIG. 1: This dartboard is a probability space with a 5-set partition. The atoms have probabilities $\frac{1}{2}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8}$. This probability distribution will be used later in Meshalkin's example on page 13.

Solve this by always dividing the remaining probability in half. 'Is it A?'; if Yes, then $V = 1$. Else: 'Is it B or C?' —if Yes, then 'Is it B?' —if No, then the dart landed in C, and $V=3$ was the number of questions. Evidently $V=3$ also for regions B, D, E. Using "log" to denote base-2 logarithm^{♥4}, the expected number of questions^{♥5} is thus

$$2: \quad \mathbb{E}(V) = \frac{1}{2} \cdot 1 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 = \sum_{j=0}^4 p_j \cdot \log\left(\frac{1}{p_j}\right) \stackrel{\text{note}}{=} 2.$$

Letting $\vec{v} := (\frac{1}{2}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8})$ be the probability vector, we can write this expectation as

$$\mathbb{E}(V) = \sum_{x \in \vec{v}} \eta(x).$$

Here, $\eta: [0, 1] \rightarrow [0, \infty)$ is the important function^{♥6}

$$3: \quad \begin{aligned} \eta(x) &:= x \cdot \log(1/x); \quad \text{so extending by continuity gives} \\ \eta(0) &= 0. \end{aligned}$$

An interpretation of " $\eta(x)$ " is the number of questions needed to winnow down to an event of probability x .

Distribution entropy

Given a probability vector \vec{v} , define its **distribution entropy** as

$$4: \quad \mathcal{H}(\vec{v}) := \sum_{x \in \vec{v}} \eta(x).$$

^{♥4}In this paper, unmarked logs will be to base-2. In entropy theory, it does not matter much what base is used, but base-2 is convenient for computing entropy for messages described in bits.

When using the natural logarithm, some people refer to the unit of information as a *nat*. In this paper, I have picked bits, rather than nats.

^{♥5}This holds when each probability p is a reciprocal power of two. For general probabilities, the "expected number of questions" interpretation holds in a weaker sense: Throw N darts independently at N copies of the dartboard. Efficiently ask Yes/No questions to determine where *all* N darts landed. Dividing by N , then sending $N \rightarrow \infty$, will be the $p \cdot \log(\frac{1}{p})$ sum of (2).

^{♥6}There does not seem to be a standard name for this function. I use η , since an uppercase η looks like an H , which is the letter that Shannon used to denote what I am calling distribution-entropy.

In this paper, I will use the term **distropy** for ‘distribution entropy’ and will reserve **entropy** for the corresponding dynamical concept, when there is a notion of *time* involved. Getting ahead of ourselves, the *entropy* of a stationary process is the asymptotic average value that its distropy decays to, as we look at larger and larger finite portions of the process.

An equi-probable vector $\vec{v} := (\frac{1}{K}, \dots, \frac{1}{K})$ evidently has $\mathcal{H}(\vec{v}) = \log(K)$. On a probability space, the “distropy of partition \mathbf{P} ”, written $\mathcal{H}(\mathbf{P})$ or $\mathcal{H}(A_1, A_2, \dots)$, shall mean the distropy of probability vector $j \mapsto \mu(A_j)$.

A (finite) partition necessarily has finite distropy. A *countable* partition can have finite distropy, e.g. $\mathcal{H}(\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \dots) = 2$. One could also have infinite distropy: Consider a piece $B \subset X$ of mass $1/2^N$. Splitting B into 2^k many equal-mass atoms gives an η -sum of $2^k \cdot \frac{k+N}{2^k 2^N}$. Setting $k = k_N := 2^N - N$ makes this η -sum equal 1; so splitting the pieces of $X = \bigsqcup_{N=1}^{\infty} B_N$, with $\mu(B_N) = \frac{1}{2^N}$, yields an ∞ -distropy partition.

The η function. Our $\eta(x) = x \cdot \log(1/x)$ function^{♥7} has vertical tangent at $x=0$, maximum at $1/e$ and, when graphed in nats, slope -1 at $x=1$.

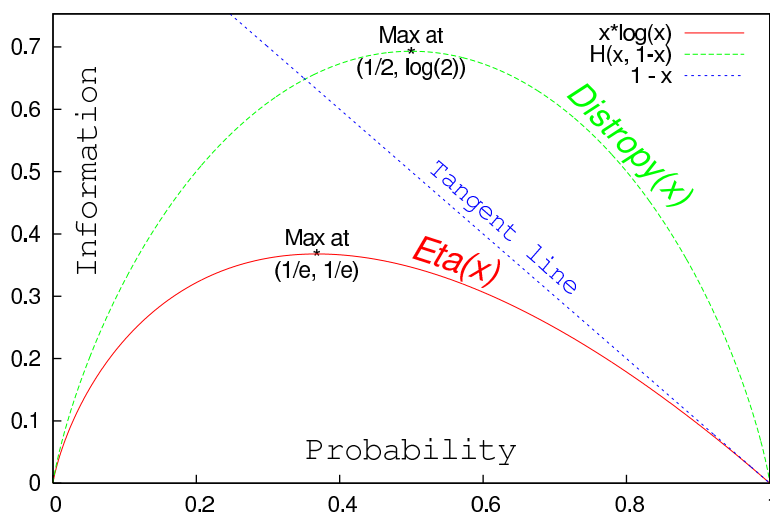


FIG. 5: Using natural log, here are the graphs of:
 $\eta(x)$ in solid red.

$\mathcal{H}(x, x^c)$ in dashed green.

$1-x$ in dotted blue.

Both $\eta(x)$ and $\mathcal{H}(x, x^c)$ are strictly convex-down. The $1-x$ line is tangent to $\eta(x)$ at $x=1$.

Consider partitions \mathbf{P} and \mathbf{Q} on the same space (X, μ) . Their *join*, written $\mathbf{P} \vee \mathbf{Q}$, has atoms $A \cap B$, for each pair $A \in \mathbf{P}$ and $B \in \mathbf{Q}$. They are *independent*, written $\mathbf{P} \perp \mathbf{Q}$, if $\mu(A \cap B) = \mu(A)\mu(B)$ for each A, B pair. We write $\mathbf{P} \geq \mathbf{Q}$, and say that “ \mathbf{P} *refines* \mathbf{Q} ”, if each \mathbf{P} -atom is a subset of some \mathbf{Q} -atom. Consequently, each \mathbf{Q} -atom is a union of \mathbf{P} -atoms.

Recall, for δ a real number, our convention that δ^c means $1 - \delta$, in analogy with $\mu(B^c)$ equaling $1 - \mu(B)$ on a probability space.

^{♥7}Curiosity: Just in this footnote we compute distropy in *nats*, that is, using natural logarithm. Given a small probability $p \in [0, 1]$ and setting $x := 1/p$, note that $\eta(p) = \frac{\log(x)}{x} \approx 1/\pi(x)$, where $\pi(x)$ denotes the number of prime numbers less-equal x . (This approximation is a weak form of the Prime Number Theorem.) Is there any actual connection between the ‘approximate distropy’ function $\mathcal{H}_\pi(\vec{p}) := \sum_{p \in \vec{p}} 1/\pi(1/p)$ and Number Theory, other than a coincidence of growth rate?

6: Distropy fact. For partitions P, Q, R on probability space (X, μ) :

a: $\mathcal{H}(P) \leq \log(\#P)$, with equality IFF P is an equi-mass partition.

b: $\mathcal{H}(Q \vee R) \leq \mathcal{H}(Q) + \mathcal{H}(R)$, with equality IFF $Q \perp R$.

c: For $\delta \in [0, \frac{1}{2}]$, the function $\delta \mapsto \mathcal{H}(\delta, \delta^c)$ is strictly increasing.

d: $R \leq P$ implies $\mathcal{H}(R) \leq \mathcal{H}(P)$, with equality IFF $R \stackrel{a.e.}{=} P$. ◇

Proof. Use the strict concavity of $\eta()$, together with Jensen's Inequality. ◇

Remark. Although we will not discuss it in this paper, most distropy statements remain true with 'partition' replaced by 'countable partition of finite distropy'. □

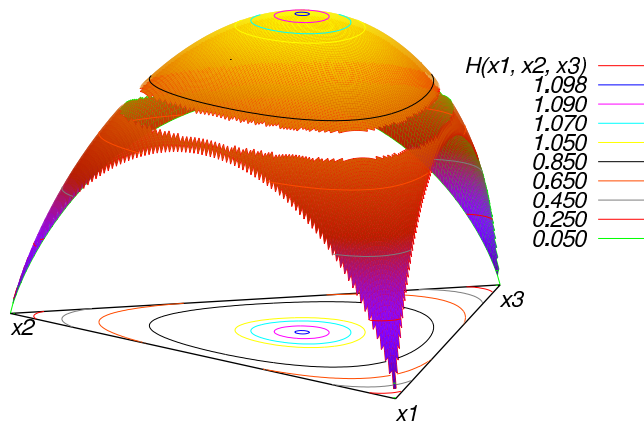


FIG. 7: Using natural log: The graph of $\mathcal{H}(x_1, x_2, x_3)$ in barycentric coordinates; a slice has been removed, between $z = 0.745$ and $z = 0.821$. The three arches are copies of the distropy curve from (5).

Binomial coefficients. The dartboard gave an example where distropy arises in a natural way. Here is a second example.

For a small $\delta > 0$, one might guess that the binomial coefficient $\binom{n}{\delta n}$ grows asymptotically (as $n \rightarrow \infty$) like $2^{\varepsilon n}$, for some small ε . But what is the correct relation between ε and δ ? Well, Stirling's formula $n! \approx [n/e]^n$ gives

$$\frac{n!}{[\delta n]! \cdot [\delta^c \cdot n]!} \approx \frac{n^n}{[\delta n]^{\delta n} \cdot [\delta^c \cdot n]^{\delta^c n}} = 1/[\delta^{\delta n} \cdot [\delta^c]^{\delta^c n}]. \quad (\text{Recall } \delta^c = 1 - \delta.)$$

Thus $\frac{1}{n} \cdot \log \binom{n}{\delta n} \approx \mathcal{H}(\delta, \delta^c)$. But by means of the above distropy inequalities, we get an inequality true for *all* n , not just asymptotically.

8: Binomial Lemma. Fix a $\delta \in [0, \frac{1}{2}]$ and let $\mathbf{H} := \mathcal{H}(\delta, \delta^c)$. Then for each $n \in \mathbb{Z}_+$:

9:
$$\sum_{j \in [0.. \delta n]} \binom{n}{j} \leq 2^{\mathbf{H}n}. \quad \diamond$$

Proof. Let $X \subset \{0, 1\}^n$ be the set of \vec{x} with $\# \{i \in [1..n] \mid x_i = 1\} \leq \delta \cdot n$. On X , let P_1, P_2, \dots be the coordinate partitions; e.g. $P_7 = (A_7, A_7^c)$, where $A_7 := \{\vec{x} \mid x_7 = 1\}$. Weighting each point by $\frac{1}{|X|}$, the uniform distribution on X , gives that $\mu(A_7) \leq \delta$. So $\mathcal{H}(P_7) \leq \mathbf{H}$, by (6c). Finally, the join $P_1 \vee \dots \vee P_n$ separates the points of X . So

$$\log(\#X) = \mathcal{H}(P_1 \vee \dots \vee P_n) \leq \mathcal{H}(P_1) + \dots + \mathcal{H}(P_n) \leq \mathbf{H}n, \quad \blacklozenge$$

making use of (6a,b). And $\#X$ equals LhS(9).

A gander at Shannon's Noisy Channel theorem

We can restate the Binomial lemma using the **Hamming metric** on $\{0, 1\}^n$,

$$\text{Dist}(\vec{x}, \vec{y}) := \# \{i \in [1..n] \mid x_i \neq y_i\}.$$

Use $\text{Bal}(\vec{x}, r)$ for the open radius- r ball centered at \vec{x} , and

$$\overline{\text{Bal}}(\vec{x}, r) := \{\vec{y} \mid \text{Dist}(\vec{x}, \vec{y}) \leq r\}$$

for the closed ball. The above lemma can be interpreted as saying that

$$9': \quad |\overline{\text{Bal}}(\vec{x}, \delta n)| \leq 2^{\mathcal{H}(\delta, \delta^c)n}, \quad \text{for each } \vec{x} \in \{0, 1\}^n.$$

10: Corollary. Fix $n \in \mathbb{Z}_+$ and $\delta \in [0, \frac{1}{2}]$, and let $\mathbf{H} := \mathcal{H}(\delta, \delta^c)$. Then there is a set $C \subset \{0, 1\}^n$, with $\#C \geq 2^{[1-\mathbf{H}]n}$, that is **strongly δn -separated**. I.e, $\text{Dist}(\vec{x}, \vec{y}) > \delta n$ for each distinct pair $\vec{x}, \vec{y} \in C$. \blacklozenge

Noisy Channel. Shannon's theorem says that a noisy channel has a **channel capacity**. Transmitting *above* this speed, there is a minimum error-rate (depending how much “above”) that no error-correcting code can fix. Conversely, one can transmit *below* –but arbitrarily close to– the channel capacity, and encode the data so as to make the error-rate less than any given ε . We use (10) to show the existence of such codes, in the simplest case where the noise[∞] is a binary independent-process (a “Bernoulli” process, in the language later in this article).

We have a channel which can pass one bit per second. Alas, there is a fixed noise-probability $\nu \in [0, \frac{1}{2})$ so that a bit in the channel is perturbed into the other value. Each perturbation is independent of all others. Let $\mathbf{H} := \mathcal{H}(\nu, \nu^c)$. The value $[1-\mathbf{H}]$ bits-per-second is the **channel capacity** of this noise-afflicted channel.

[∞]The noise-process is assumed to be *independent* of the signal-process. In contrast, when the perturbation is highly dependent on the signal, then it is sometimes called **distortion**.

Encoding/Decoding. Encode using an “ k, n -block-code”; an injective map $F: \{0, 1\}^k \rightarrow \{0, 1\}^n$. The source text is split into consecutive k -bit blocks. A block $\vec{x} \in \{0, 1\}^k$ is encoded to $F(\vec{x}) \in \{0, 1\}^n$ and then sent through the channel, where it comes out perturbed to $\vec{\alpha} \in \{0, 1\}^n$. The *transmission rate* is thus k/n bits-per-second.

For this example, we fix a radius $r > 0$ to determine the decoding map,

$$D_r: \{0, 1\}^n \rightarrow \{Oops\} \sqcup \{0, 1\}^k.$$

We set $D_r(\vec{\alpha})$ to \vec{z} if there is a *unique* \vec{z} with $F(\vec{z}) \in \overline{\text{Bal}}(\vec{\alpha}, r)$; else, set $D_r(\vec{\alpha}) := \text{Oops}$.

We can think of the noise as a $\{0, 1\}$ -independent-process, with $\text{Prob}(1) = \nu$, which is added mod-2 to the signal-process. Suppose we can arrange that the set $\{F(\vec{x}) \mid \vec{x} \in \{0, 1\}^k\}$ of codewords, is a strongly r -separated-set. Then

11: *The probability that a block is mis-decoded is the probability, flipping a ν -coin n times, that we get more than r many Heads.*

12: Theorem (Shannon). Fix a noise-probability $\nu \in [0, \frac{1}{2})$ and let $\mathbf{H} := \mathcal{H}(\nu, \nu^c)$. Consider a rate $R < [1 - \mathbf{H}]$ and an $\varepsilon > 0$. Then $\forall_{\text{large } n}$ there exists a k and a code $F: \{0, 1\}^k \rightarrow \{0, 1\}^n$ so that: The F -code transmits bits at faster than R bits-per-second, and with error-rate $< \varepsilon$. ♦

Proof. Let $\mathbf{H}' := \mathcal{H}(\delta, \delta^c)$, where $\delta > \nu$ was chosen so small that

$$13: \quad \delta < \frac{1}{2} \quad \text{and} \quad 1 - \mathbf{H}' > R.$$

Pick a large n for which

$$14: \quad \frac{k}{n} > R, \quad \text{where } k := \lfloor [1 - \mathbf{H}'] \cdot n \rfloor.$$

By (10), there is a strongly δn -separated-set $C \subset \{0, 1\}^n$ with $\#C \geq 2^{[1 - \mathbf{H}']n}$. So C is big enough to permit an injection $F: \{0, 1\}^k \hookrightarrow C$. Courtesy (11), the probability of a decoding error is that of getting more than δn many Heads in flipping a ν -coin n times. Since $\delta > \nu$, the Weak Law of Large Numbers guarantees –once n is large enough– that this probability is less than the given ε . ♦

The information function

We use $\mathbf{P} = (A_1, \dots)$, $\mathbf{Q} = (B_1, \dots)$, $\mathbf{R} = (C_1, \dots)$ for partitions, and \mathcal{F}, \mathcal{G} for fields.

With \mathfrak{C} a (finite or infinite) family of subfields of \mathcal{X} , their *join* $\bigvee_{\mathcal{G} \in \mathfrak{C}} \mathcal{G}$ is the smallest field \mathcal{F} such that $\mathcal{G} \subset \mathcal{F}$, for each $\mathcal{G} \in \mathfrak{C}$. A partition \mathbf{Q} can be interpreted also as a field; namely, the field of unions of its atoms. A join of denumerably many partitions will be interpreted as a field, but a join of *finitely* many, $\mathbf{P}_1 \vee \dots \vee \mathbf{P}_N$, will be viewed as a partition *or* as a field, depending on context.

For an $A \subset X$, use $\mathbf{1}_A: X \rightarrow \{0, 1\}$ for its *indicator function*; $\mathbf{1}_A(x) = 1$ IFF $x \in A$. The *information function* of partition P , a map $I_P: X \rightarrow [0, \infty)$, is

$$15: \quad I_P(\cdot) := \sum_{A \in P} \log\left(\frac{1}{\mu(A)}\right) \cdot \mathbf{1}_A(\cdot).$$

It has been defined so that its expectation is the distropy of P .

$$\mathbb{E}(I_P) = \int_X I_P(\cdot) d\mu = \mathcal{H}(P).$$

With respect to a subfield \mathcal{F} , let $\mu(A \mid \mathcal{F})$ be the *conditional probability* function; that is, the conditional expectation $\mathbb{E}(\mathbf{1}_A \mid \mathcal{F})$. This engenders the *conditional information function*,

$$16: \quad \begin{aligned} I_{P \mid \mathcal{F}}(x) &:= \sum_{A \in P} \log\left(\frac{1}{\mu(A \mid \mathcal{F})(x)}\right) \cdot \mathbf{1}_A(x). \quad \text{Its integral} \\ \mathcal{H}(P \mid \mathcal{F}) &:= \int I_{P \mid \mathcal{F}} d\mu, \quad \text{is the } \mathbf{conditional\ distropy} \text{ of } P \text{ on } \mathcal{F}. \end{aligned}$$

Conditioning on a positive-mass set B , let $P|B$ be the probability vector $A \mapsto \frac{\mu(A \cap B)}{\mu(B)}$. Conditional distropy, when conditioning on a partition, equals

$$17: \quad \mathcal{H}(P \mid Q) = \sum_{B \in Q} \mathcal{H}(P|B) \cdot \mu(B) = \sum_{A \in P, B \in Q} \log\left(\frac{1}{\mu(A \cap B)/\mu(B)}\right) \cdot \mu(A \cap B).$$

Write $\mathcal{G}_j \nearrow \mathcal{F}$ to indicate that fields $\mathcal{G}_1 \subset \mathcal{G}_2 \subset \dots$ are nested, and that $\text{Fld}(\bigcup_1^\infty \mathcal{G}_j) = \mathcal{F}$, a.e. The Martingale Convergence Theorem, [20, P. 103], gives (c), below.

18: Conditional-distropy fact. Consider partitions P, Q, R and fields \mathcal{F} and \mathcal{G}_j . Then

a: $0 \leq \mathcal{H}(P \mid \mathcal{F}) \leq \mathcal{H}(P)$, with equality IFF $P \stackrel{\text{a.e.}}{\subset} \mathcal{F}$, respectively, $P \perp \mathcal{F}$.

b: $\mathcal{H}(Q \vee R \mid \mathcal{F}) \leq \mathcal{H}(Q \mid \mathcal{F}) + \mathcal{H}(R \mid \mathcal{F})$.

c: Suppose $\mathcal{G}_j \nearrow \mathcal{F}$. Then $\mathcal{H}(P \mid \mathcal{G}_j) \searrow \mathcal{H}(P \mid \mathcal{F})$.

d: $\mathcal{H}(Q \vee R) = \mathcal{H}(Q \mid R) + \mathcal{H}(R)$.

d': $\mathcal{H}(Q \vee R_1 \mid R_0) = \mathcal{H}(Q \mid R_1 \vee R_0) + \mathcal{H}(R_1 \mid R_0)$. ◇

Imagining our dartboard (1) divided by superimposed partitions Q and R , equality (18d) can interpreted as saying: ‘You can efficiently discover where the dart landed in both partitions, by first asking efficient questions about R , then –based on where you landed in R – asking intelligent questions about Q .’

Entropy of a Process

Consider an transformation $(T : X, \mu)$ and partition $P = (A_1, A_2, \dots)$. Each “time” n determines a partition $P_n := T^n P$, whose j^{th} -atom is $T^{-n}(A_j)$. The **process** T, P refers to how T acts on the subfield $\bigvee_0^\infty P_n \subset \mathcal{X}$. (An alternative view of a process is as a stationary sequence V_0, V_1, \dots of random variables $V_n : X \rightarrow \mathbb{Z}_+$, where $V_n(x) := j$ because x is in the j^{th} -atom of P_n .)

Write $\mathcal{E}(T, P)$ or $\mathcal{E}^T(P)$ for the “**entropy** of the T, P process”. It is the limit of the **conditional-distropy-numbers**

$$c_n := \mathcal{H}(P_0 \mid P_1 \vee P_2 \vee \dots \vee P_{n-1}).$$

This limit exists since $\mathcal{H}(P) = c_1 \geq c_2 \geq \dots \geq 0$.

Define the **average-distropy-number** $\frac{1}{n}h_n$, where

$$h_n := \mathcal{H}(P_0 \vee P_1 \vee \dots \vee P_{n-1}).$$

Certainly $h_n = c_n + \mathcal{H}(P_1 \vee \dots \vee P_{n-1}) = c_n + h_{n-1}$, since T is measure preserving. Induction gives $h_n = \sum_{j=1}^n c_j$. So the Cesàro averages $\frac{1}{n}h_n$ converge to the entropy.

19: Theorem. *The entropy of process $(T, P : X, \mathcal{X}, \mu)$ equals*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \mathcal{H}(P_0 \vee \dots \vee P_{n-1}) = \lim_{n \rightarrow \infty} \mathcal{H}(P_0 \mid \bigvee_1^n P_j) = \mathcal{H}(P_0 \mid \bigvee_1^\infty P_j).$$

Both limits are non-increasing. The entropy $\mathcal{E}^T(P) \geq 0$, with equality IFF $P \stackrel{a.e.}{\subset} \bigvee_1^\infty P_j$. And $\mathcal{E}^T(P) \leq \mathcal{H}(P)$, with equality IFF T, P is an independent process. \diamond

Generators. We henceforth only discuss *invertible* mpts, that is, when T^{-1} is itself an mpt. Viewing the atoms of P as “letters”, then, each $x \in X$ has a **T, P -name** $\dots x_{-1} x_0 x_1 x_2 \dots$, where x_n is $P(T^n(x))$, the P -letter owning $T^n(x)$.

A partition P **generates** (the whole field) under $(T : X, \mu)$, if $\bigvee_{-\infty}^\infty T^n P =_\mu \mathcal{X}$. It turns out⁹ that P generates IFF P **separates points**. That is, after deleting a (T -invariant) nullset, distinct points of X have distinct T, P -names.

A finite set $[1 \dots L]$ of integers, our **alphabet**, yields the **shift space** $X := [1 \dots L]^\mathbb{Z}$ of doubly-infinite sequences $x = (\dots x_{-1} x_0 x_1 \dots)$. The **shift** $T : X \rightarrow X$ acts on X by

$$T(x) := [n \mapsto x_{n+1}].$$

Automatically, then, the **time-zero partition** P separates points, under the action of the shift. This L -atom partition has $P\langle x \rangle = P\langle y \rangle$ IFF $x_0 = y_0$. So no matter what shift-invariant measure is put on X , the time-zero partition will generate under the action of T .

⁹I am now at liberty to reveal that our X has always been a **Lebesgue space**, that is, measure-isomorphic to an interval of \mathbb{R} together with countably many point-atoms (points with positive mass). The equivalence of *generating* and *separating* is a technical theorem, due to Rokhlin.

Assuming μ to be Lebesgue is not much of a limitation. For instance, if μ is a finite measure on any Polish space, then μ extends to a Lebesgue measure on the μ -completion of the Borel sets. To not mince words: All spaces are Lebesgue spaces unless you are actively *looking* for trouble...

Time reversibility. A transformation need not be isomorphic to its inverse. Nonetheless, the average-distropy-numbers show that $\mathcal{E}(T^{-1}, \mathbf{P}) = \mathcal{E}(T, \mathbf{P})$; although this is not obvious from the conditioning-definition of entropy. Alternatively,

$$\begin{aligned} 20: \quad \mathcal{H}(\mathbf{P}_0 \mid \bigvee_1^n \mathbf{P}_j) &= \mathcal{H}(\mathbf{P}_0 \vee \dots \vee \mathbf{P}_n) - \mathcal{H}(\mathbf{P}_1 \vee \dots \vee \mathbf{P}_n) \\ &= \mathcal{H}(\mathbf{P}_{-n} \vee \dots \vee \mathbf{P}_0) - \mathcal{H}(\mathbf{P}_{-n} \vee \dots \vee \mathbf{P}_1) = \mathcal{H}(\mathbf{P}_0 \mid \bigvee_{-1}^{-n} \mathbf{P}_j). \end{aligned} \quad \square$$

Bernoulli processes. A probability vector $\vec{v} := (v_1, \dots, v_L)$ can be viewed as a measure on alphabet $[1 \dots L]$. Let $\mu_{\vec{v}}$ be the resulting product measure on $X := [1 \dots L]^{\mathbb{Z}}$, with T the shift on X and \mathbf{P} the time-zero partition. The independent process $(T, \mathbf{P} : X, \mu_{\vec{v}})$ is called, by ergodic theorists, a **Bernoulli process**. Not necessarily consistently, we tend to refer to the underlying transformation as a **Bernoulli shift**.

The $(\frac{1}{2}, \frac{1}{2})$ -Bernoulli and the $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ -Bernoulli have different process-entropies, but perhaps their underlying transformations are isomorphic? Prior to the Kolmogorov-Sinai definition of entropy^{♥10} of a transformation, this question remained unanswered.

Entropy of a Transformation

The Kolmogorov-Sinai definition of the entropy of an mpt is

$$\mathcal{E}(T) := \sup\{\mathcal{E}^T(\mathbf{Q}) \mid \mathbf{Q} \text{ a partition on } X\}.$$

Certainly entropy is an isomorphism invariant—but is it useful? After all, the supremum of *distropies* of partitions is always infinite (on non-atomic spaces) and one might fear that the same holds for entropies. The key observation (restated in (23c) and proved below) was this, from [4, Kol 1958] and [5, Sinai 1959].

21: Kolmogorov-Sinai theorem. *If \mathbf{P} generates under T , then $\mathcal{E}(T) = \mathcal{E}(T, \mathbf{P})$.* ♦

Thereupon the $(\frac{1}{2}, \frac{1}{2})$ and $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ Bernoulli-shifts are *not* isomorphic, since their respective entropies are $\log(2) \neq \log(3)$.

Wolfgang Krieger later proved a converse to the Kolmogorov-Sinai theorem.

22: Krieger Generator Thm (1970). *Suppose T ergodic. If $\mathcal{E}(T) < \infty$, then T has a generating partition. Indeed, letting K be the smallest integer $K > \mathcal{E}(T)$, there is a K -atom generator.*^{♥11} ♦

Proof. See Rudolph [21], or [[de la Rue, §5.1]], where Krieger's theorem is stated in terms of joinings. ♦

^{♥10}This is sometimes called *measure(-theoretic) entropy* or (perhaps unfortunately) *metric entropy*, to distinguish it from topological entropy. Tools known prior to entropy, such as *spectral* properties, did *not* distinguish the two Bernoulli-shifts; see [[Lemanczyk]] for the definitions.

^{♥11}It is an easier result, undoubtedly known much earlier, that every ergodic T has a *countable* generating partition—possibly of ∞ -distropy.

Entropy is continuous. Given *ordered* partitions $Q = (B_1, \dots)$ and $Q' = (B'_1, \dots)$, extend the shorter by null-atoms until $|Q| = |Q'|$. Let $Fat := \bigsqcup_j [B_j \cap B'_j]$; this set should have mass close to 1 if Q and Q' are almost the same partition. Define a new partition

$$Q \triangle Q' := \{Fat\} \sqcup \{B_i \cap B'_j \mid \text{with } i \neq j\}.$$

(In other words, take $Q \vee Q'$ and coalesce, into a single atom, all the $B_k \cap B'_k$ sets.) Topologize the space of partitions by saying^{♥12} that $Q^{(L)} \rightarrow Q$ when $\mathcal{H}(Q \triangle Q^{(L)}) \rightarrow 0$. Then (23b) says that process-entropy varies continuously with varying the partition.

23: Lemma. Fix a mpt $(T : X, \mu)$. For partitions P, Q, Q' , define $R := Q \triangle Q'$ and let $\delta := \mathcal{H}(R)$. Then

a: $|\mathcal{H}(Q) - \mathcal{H}(Q')| \leq \delta$. (Distropy varies continuously with the partition.)

b: $|\mathcal{E}^T(Q) - \mathcal{E}^T(Q')| \leq \delta$. (Process-entropy varies continuously with the partition.)

c: For all partitions $Q \subset \text{Fld}(T, P)$: $\mathcal{E}^T(Q) \leq \mathcal{E}^T(P)$. ♦

Proof of (a). Evidently $Q' \vee R = Q' \vee Q = Q \vee R$. So $\mathcal{H}(Q') \leq \mathcal{H}(Q \vee R) \leq \mathcal{H}(Q) + \delta$. ♦

Proof of (b). As above, $\mathcal{H}(\bigvee_1^N Q'_j) \leq \mathcal{H}(\bigvee_1^N Q_j) + \mathcal{H}(\bigvee_1^N R_j)$. Sending $N \rightarrow \infty$ gives $\mathcal{E}^T(Q') \leq \mathcal{E}^T(Q) + \mathcal{E}^T(R)$. Finally, $\mathcal{E}^T(R) \leq \mathcal{H}(R)$ and so $\mathcal{E}^T(Q') \leq \mathcal{E}^T(Q) + \delta$. ♦

Proof of (c). Let $K := |Q|$. Then there is a sequence of K -set partitions $Q^{(L)} \rightarrow Q$ with $Q^{(L)} \leq \bigvee_{-L}^L P_\ell$. By above, $\mathcal{E}^T(Q^{(L)}) \rightarrow \mathcal{E}^T(Q)$, so showing that $\mathcal{E}^T(\bigvee_{-L}^L P_\ell) \stackrel{?}{\leq} \mathcal{E}^T(P)$ will suffice. Note that

$$h_N := \mathcal{H}\left(\bigvee_{n=0}^{N-1} T^n\left(\bigvee_{-L}^L P_\ell\right)\right) = \mathcal{H}\left(\bigvee_{j=-L}^{N-1+L} P_j\right).$$

So $\frac{1}{N}h_N \leq \frac{1}{N}\mathcal{H}\left(\bigvee_0^{N-1} P_j\right) + \frac{1}{N} \cdot 2L \cdot \mathcal{H}(P)$. Now send $N \rightarrow \infty$. ♦

Entropy is not continuous. The most common topology placed on the space, Ω , of mpts is the *coarse topology*^{♥13} that Halmos discusses in his “little red book”, [14].

The Rokhlin lemma [21, P. 33] implies that the isomorphism-class of *each* ergodic mpt is *dense* in Ω , (e.g, see [14, P. 77]) disclosing that the $S \mapsto \mathcal{E}(S)$ map is exorbitantly discontinuous.

Indeed, the failure happens already for process-entropy with respect to a fixed partition. A Bernoulli process T, P has positive entropy. Take mpts $S_n \rightarrow T$, each isomorphic to an irrational rotation. Then each $\mathcal{E}(S_n, P)$ is zero, as shown in the later section on *Determinism and Zero-entropy*.

^{♥12}On the set of ordered K -set partitions (with K fixed) this convergence is the same as: $Q^{(L)} \rightarrow Q$ when $\mu(Fat(Q^{(L)}, Q)) \rightarrow 1$.

An alternative approach is the **Rokhlin metric**, $\text{Dist}(P, Q) := \mathcal{H}(P \mid Q) + \mathcal{H}(Q \mid P)$, which has the advantage of working for *unordered* partitions.

^{♥13}I.e, $S_n \rightarrow T$ IFF $\forall A \in \mathcal{X} : \mu(S_n^{-1}(A) \triangle T^{-1}(A)) \rightarrow 0$; this is a metric-topology, since our probability space is countably generated. This can be restated in terms of the unitary operator U_T on $\mathbb{L}^2(\mu)$, where $U_T(f) := f \circ T$. Namely, $S_n \rightarrow T$ in the coarse topology IFF $U_{S_n} \rightarrow U_T$ in the strong operator topology.

Further results. When \mathcal{F} is a T -invariant subfield, agree to use $T \upharpoonright_{\mathcal{F}}$ for “ T restricted to \mathcal{F} ”, which is a factor (see Glossary) of T . Transformations T and S are **weakly isomorphic** if each is isomorphic to a factor of the other.

The foregoing entropy tools make short shrift of the following.

24: Entropy lemma. Consider T -invariant subfields \mathcal{G}_j and \mathcal{F} .

a: Suppose $\mathcal{G}_j \nearrow \mathcal{F}$. Then $\mathcal{E}(T \upharpoonright_{\mathcal{G}_j}) \nearrow \mathcal{E}(T \upharpoonright_{\mathcal{F}})$. In particular, $\mathcal{G} \subset \mathcal{F}$ implies that $\mathcal{E}(T \upharpoonright_{\mathcal{G}}) \leq \mathcal{E}(T \upharpoonright_{\mathcal{F}})$, so entropy is an invariant of weak-isomorphism.

b: $\mathcal{E}(T \upharpoonright_{\mathcal{G}_1 \vee \mathcal{G}_2 \vee \dots}) \leq \sum_j \mathcal{E}(T \upharpoonright_{\mathcal{G}_j})$. And $\mathcal{E}(T, Q_1 \vee Q_2 \vee \dots) \leq \sum_j \mathcal{E}(T, Q_j)$.

c: For mpts $(S_j : Y_j, \nu_j)$: $\mathcal{E}(S_1 \times S_2 \times \dots) = \sum_j \mathcal{E}(S_j)$.

d: $\mathcal{E}(T^{-1}) = \mathcal{E}(T)$. More generally, $\mathcal{E}(T^n) = |n| \cdot \mathcal{E}(T)$. ◇

E.g: Meshalkin’s map. In the wake of Kolmogorov’s 1958 entropy paper, for two Bernoulli-shifts to be isomorphic one now knew that they had to have equal entropies. Meshalkin provided the first non-trivial example [44], in 1959.

Let $S:Y \curvearrowright$ be the Bernoulli-shift over the “letter” alphabet $\{E, D, P, N\}$, with probability distribution $(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$. The letters E, D, P, N stand for *Even, oDd, Positive, Negative*, and will be used to describe the code (isomorphism) between the processes.

Use $T:X \curvearrowright$ for the Bernoulli-shift over “digit” alphabet $\{0, +1, -1, +2, -2\}$, with probability distribution $(\frac{1}{2}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8})$. Both distributions, $(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$ and $(\frac{1}{2}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8})$, have dis-
entropy $\log(4)$.

After deleting invariant nullsets from X and Y , we will construct a measure-preserving isomorphism $\psi: X \rightarrow Y$ so that $T \circ \psi = \psi \circ S$.

The Code. In X , consider this point x :

... 0 0 0 -1 0 0 +1 +2 -1 +1 0 ...

Regard each 0 as a left-parenthesis, and each non-zero as a right-parenthesis. Link them according to the legal way of matching parentheses, as shown in the top row, below:

0	0	0	-1	0	0	+1	+2	-1	+1	0
P	N	N	D	P	P	D	E	D	D	?

The leftmost 0 is linked to the rightmost +1, as indicated by the longest-overbar. The left/right-parentheses form a $(\frac{1}{2}, \frac{1}{2})$ -random-walk. Since this random walk is recurrent, we know that every position in x will be linked (except for a nullset of points x).

Below each **0**, write “P” or “N” as the **0** is linked to a *positive* or *negative* digit. And below the other digits, write “E” or “D” as the digit is *even* or *odd*. So the upper name in X is mapped to the lower name, a point $y \in Y$.

This map $\psi: X \rightarrow Y$ carries the upstairs $(\frac{1}{2}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8})$ distribution to $(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$, downstairs. It takes some arguing to show that independence is preserved.

The inverse map, ψ^{-1} , views D and E as right-parentheses, and P and N as left. Above D, write the odd digit **+1** or **-1**, as this D is linked to Positive or Negative. \square

E.g: Markov Shifts. A Bernoulli process T, P has independence, $P_{(-\infty..0]} \perp P_1$, whereas a *Markov process* is a bit less aloof:

The infinite Past $P_{(-\infty..0]}$ doesn't provide any more information about Tomorrow than Today did.

That is, the conditional distribution $P_1 | P_{(-\infty..0]}$ equals $P_1 | P_0$. Equivalently,

$$25: \quad \mathcal{H}(P_1 | P_0) = \mathcal{H}(P_1 | P_{(-\infty..0]}) \stackrel{\text{note}}{=} \mathcal{E}(T, P).$$

The simplest non-trivial Markov process $(T, P : X, \mu)$ is over a two-letter alphabet $\{a, b\}$, and has transition graph (26), for some choice of transition probabilities s and c . The

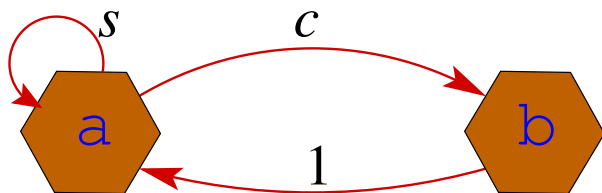


FIG. 26: Call the transition probabilities $s := \text{Prob}(a \rightarrow a)$ for stay, and $c := \text{Prob}(a \rightarrow b)$ for change. These are non-negative reals, and $s+c = 1$.

graph's Markov matrix is

$$M = [m_{i,j}]_{i,j} = \begin{bmatrix} s & c \\ 1 & 0 \end{bmatrix}, \quad \text{where } c = 1 - s, \text{ and } m_{i,j} \text{ denotes the probability of going from state } i \text{ to state } j.$$

If Today's distribution on the two states is the probability-vector $\vec{v} := [p_a \ p_b]$, then Tomorrow's is the product $\vec{v} \cdot M$. So a *stationary* process needs $\vec{v}M = \vec{v}$. This equation has the unique solution $p_a = \frac{1}{1+c}$ and $p_b = \frac{c}{1+c}$. An example of computing the probability of a word (or cylinder set; see [Petersen, 5.1]) in the process, is

$$\mu_s(\text{baaaba}) = p_b \cdot m_{ba} \cdot m_{aa} \cdot m_{aa} \cdot m_{ab} \cdot m_{ba} = \frac{c}{1+c} \cdot 1 \cdot s \cdot s \cdot c \cdot 1.$$

The subscript on μ_s indicates the dependence on the transition probabilities; let's also mark the mpt and call it T_s . Using (25), the entropy of our Markov map is

$$27: \quad \mathcal{E}(T_s) = p_a \cdot \mathcal{H}(s, c) + p_b \cdot \mathcal{H}(1, 0) = \frac{1}{1+c} \cdot [s \log(s) + c \log(c)].$$

\square

Determinism and Zero-entropy

Irrational rotations have zero-entropy; let's reveal this in two different ways.

Equip $X := [0, 1)$ with “length” (Lebesgue) measure and wrap it into a circle. With “ \oplus ” denoting addition mod-1, have $T: X \rightarrow X$ be the rotation $T(x) := x \oplus \alpha$, where the rotation number α is irrational. Pick distinct points $y_0, z_0 \in X$, and let P be the partition whose two atoms are the intervals $[y_0, z_0)$ and $[z_0, y_0)$, wrapping around the circle.

The T -orbit of each point x is dense^{♥14} in X . In particular, y_0 has dense orbit, so P separates points –hence generates– under T . Our goal, thus, is $\mathcal{E}(T, P) \stackrel{?}{=} 0$.

E.g. Rotations are deterministic. The forward T -orbit of each point is dense. This is true for y_0 , and so the *backward* T, P -name of each x actually tells us which point x is. I.e., $P \subset \bigvee_{n=-\infty}^{-1} T^n P$, which is our definition of “process T, P is **deterministic**”. Our P being finite, this determinism implies that $\mathcal{E}(T, P)$ is zero, by (19).

Counting names in a rotation. The $P_0 \vee \dots \vee P_{n-1}$ partition places n translates of points y_0 and of z_0 , cutting the circle into at most $2n$ intervals. Thus $\mathcal{H}(P_0 \vee \dots \vee P_{n-1}) \leq \log(2n)$. And $\frac{1}{n} \log(2n) \rightarrow 0$.

Alternatively, the below SMB-theorem implies, for an ergodic process T, P , that the number of length- n names is approximately $2^{\mathcal{E}(T, P)n}$; this, after discarding small mass from the space. But the growth of $n \mapsto 2n$ is sub-exponential and so, for our rotation, $\mathcal{E}(T, P)$ must be zero. \square

28: Shannon-McMillan-Breiman^{♥15} Theorem (SMB-Thm). Set $E := \mathcal{E}(T, P)$, where tuple $(T, P : X, \mu)$ is an ergodic process. Then the average information function

$$28a: \quad \frac{1}{n} \cdot \mathcal{I}_{P_{[0..n)}}(x) \xrightarrow{n \rightarrow \infty} E, \quad \text{for a.e } x \in X.$$

The functions $f_n := \mathcal{I}_{P_{[0..n)}}$ converge to the constant function E both in the \mathbb{L}^1 -norm and in probability. \diamond

Consequences. Recall that $P_{[0..n)}$ means $P_0 \vee P_1 \vee \dots \vee P_{n-1}$, where $P_j := T^j P$. As usual, $P_{[0..n)} \langle x \rangle$ denotes the $P_{[0..n)}$ -atom owning x .

Having deleted a nullset, we can restate (28a) to now say that $\forall \varepsilon, \forall x, \forall_{\text{large } n}$:

$$28b: \quad 1/2^{[E+\varepsilon]n} \leq \mu(P_{[0..n)} \langle x \rangle) \leq 1/2^{[E-\varepsilon]n}.$$

This has the following consequence. Fixing a number $\delta > 0$, we consider any set with $\mu(B) \geq \delta$ and count the number of n -names of points in B . The SMB-Thm implies

$$28c: \quad \forall \varepsilon, \forall_{\text{large } n}, \forall B \stackrel{\mu}{\geq} \delta : \quad |\{n\text{-names in } B\}| \geq 2^{[E-\varepsilon]n}.$$

^{♥14}Fix an $\varepsilon > 0$ and an $N > 1/\varepsilon$. Points $x, T(x), \dots, T^N(x)$ have some two at distance less than $\frac{1}{N}$; say, $\text{Dist}(T^i(x), T^j(x)) < \varepsilon$, for some $0 \leq i < j \leq N$. Since T is an isometry, $\varepsilon > \text{Dist}(x, T^k(x)) > 0$, where $k := j - i$. So the T^k -orbit of x is ε -dense.

^{♥15}In engineering circles, this is called the Almost-everywhere equi-partition theorem.

E.g: Rank-1 has zero-entropy. There are several equivalent definitions for “rank-1 transformation”, several of which are discussed in the introduction of [28]. (See [13, Chap. 6] and [47] and [27] for examples of stacking constructions.)

A **rank-1 transformation** $(T : X, \mu)$ admits a generating partition \mathbf{P} and a sequence of Rokhlin stacks $S_n \subset X$, with heights going to ∞ , and with $\mu(S_n) \rightarrow 1$. Moreover, each of these Rokhlin stacks is \mathbf{P} -monochromatic, that is, each level of the stack lies entirely in some atom of \mathbf{P} .

Taking a stack of some height $2n$, let $B=B_n$ be the union of the bottom n levels of the stack. There are at most n many length- n names starting in B_n , by monochromaticity. Finally, $\mu(B_n)$ is almost $\frac{1}{2}$, so is certainly larger than $\delta := \frac{1}{3}$. Thus (28c) shows that our rank-1 T has zero entropy. \square

Cautions on determinism’s relation to zero-entropy. A finite-valued process T, \mathbf{P} has zero-entropy iff $\mathbf{P} \subset \bigvee_{-\infty}^{-1} \mathbf{P}_j$. Iterating gives

$$\bigvee_0^{\infty} \mathbf{P}_j \subset \bigvee_{-\infty}^{-1} \mathbf{P}_j, \quad \text{i.e, the future is measurable with respect to the past.}$$

This was the case with the rotation, where a point’s past uniquely identified the point, thus telling us its future.

While determinism and zero-entropy mean the same thing for finite-valued processes, this fails catastrophically for real-valued (i.e, continuum-valued) processes, as shown by an example of the author’s. A stationary real-valued process $\mathbf{V} = \dots V_{-1} V_0 V_1 V_2 \dots$ is constructed in [39] which is simultaneously

strongly deterministic : The two values V_0, V_1 determine all of \mathbf{V} , future and past.

and **non-consecutively independent**. This latter means that for each bi-infinite increasing integer sequence $\{n_j\}_{j=-\infty}^{\infty}$ with no consecutive pair (always $1 + n_j < n_{j+1}$), then the list of random variables $\dots V_{n_{-1}} V_{n_0} V_{n_1} V_{n_2} \dots$ is an independent process.

Restricting the random variables to be *countably*-valued, how much of the example survives? Joint work with Kalikow, [40], produced a countably-valued stationary \mathbf{V} which is non-consecutively independent as well as deterministic. (Strong determinism is ruled out, due to cardinality considerations.) A side-effect of the construction is that \mathbf{V} ’s time-reversal $n \mapsto V_{-n}$ is **not** deterministic.

Pinsker Field. Define a collection of sets (the script \mathcal{Z} is for “zero”)

$$29: \quad \mathcal{Z} = \mathcal{Z}_T = \mathcal{Z}(T) := \{D \in \mathcal{X} \mid \mathcal{E}(T, (D, D^c)) = 0\}.$$

Courtesy (24b), \mathcal{Z} is a T -invariant *field*, and

$$29': \quad \forall Q \subset \mathcal{Z} : \quad \mathcal{E}(T, Q) = 0.$$

The **Pinsker field** of $T^{\heartsuit 16}$ is this \mathcal{Z} . It is maximal with respect to (29'). Unsurprisingly, the **Pinsker factor** $T \upharpoonright_{\mathcal{Z}}$ has zero entropy, that is, $\mathcal{E}(T \upharpoonright_{\mathcal{Z}}) = 0$.

The asymptotic past of the T, P process is called its **tail field**, where

$$\text{Tail}(T, P) := \bigcap_{L=1}^{\infty} \bigvee_{j=-\infty}^{-L} P_j.$$

30: Theorem (Pinsker). Suppose P is a generating partition for an ergodic T . Then $\text{Tail}(T, P)$ equals \mathcal{Z}_T . In particular, all generating partitions for T have the same tail field. And the **future field** of T, P , which is $\mathcal{Z}(T^{-1}) \stackrel{\text{note}}{=} \mathcal{Z}(T)$, equals its tail field. \diamond

Ornstein theory

In 1970, Don Ornstein solved the long-standing problem of showing that entropy was a complete isomorphism-invariant of Bernoulli transformations; that is, that two independent processes with same entropy necessarily have the same underlying transformation. (Earlier, Sinai had shown that two such Bernoulli maps were *weakly isomorphic*, that is, each isomorphic to a factor of the other.)

Ornstein introduced the notion of a process being *finitely determined*, see [46] for a definition, proved that a transformation T was Bernoulli IFF it had a finitely-determined generator IFF every partition was finitely-determined with respect to T , and showed that entropy completely classified the finitely-determined processes upto isomorphism.

The Pinsker-Field and K-automorphisms

Said differently, the zero-entropy trns are those whose Pinsker-field is everything. TBW ●●●

¹⁶Traditionally, this called the *Pinsker algebra* where, in this context, “algebra” is understood to mean “ σ -algebra”.

Topological Entropy

Adler, Konheim and McAndrew, in 1965, published the first definition of *topological entropy* in the eponymous article [32]. Here, $T:X \rightarrow X$ is a continuous self-map of a compact topological space. The role of atoms is played by open sets. Instead of a finite partition, one uses a finite^{♥17} **open-cover** $\mathcal{V} = \{U_j\}_{j=1}^L$, i.e each **patch** U_j is open, and their union $\bigcup(\mathcal{V}) = X$. (Henceforth, ‘cover’ means “open cover”.)

Let $\text{Card}(\mathcal{V})$ be the minimum cardinality over all subcovers.

$$\text{Card}(\mathcal{V}) := \text{Min} \{ \# \mathcal{V}' \mid \mathcal{V}' \subset \mathcal{V} \text{ and } \bigcup(\mathcal{V}') = X \}, \text{ and let}$$

$$\mathcal{H}(\mathcal{V}) = \mathcal{H}_{\text{top}}(\mathcal{V}) := \log(\text{Card}(\mathcal{V})).$$

Analogous to the definitions for partitions, we prescribe

$$\begin{aligned} \mathcal{V} \vee \mathcal{W} &:= \{V \cap W \mid V \in \mathcal{V} \text{ and } W \in \mathcal{W}\}; \\ T\mathcal{V} &:= \{T^{-1}(U) \mid U \in \mathcal{V}\} \text{ and } \mathcal{V}_{[0..n)} := \mathcal{V}_0 \vee \mathcal{V}_1 \vee \dots \vee \mathcal{V}_{n-1}; \\ \mathcal{W} \geq \mathcal{V} &, \text{ if each } \mathcal{W}\text{-patch is a subset of some } \mathcal{V}\text{-patch.} \end{aligned}$$

The T, \mathcal{V} -**entropy** is

$$31: \quad \mathcal{E}^T(\mathcal{V}) = \mathcal{E}(T, \mathcal{V}) = \mathcal{E}_{\text{top}}(T, \mathcal{V}) := \limsup_{n \rightarrow \infty} \frac{1}{n} \cdot \mathcal{H}_{\text{top}}(\mathcal{V}_{[0..n)}).$$

And the **topological entropy** of T is

$$32: \quad \mathcal{E}_{\text{top}}(T) := \sup_{\mathcal{V}} \mathcal{E}_{\text{top}}(T, \mathcal{V}), \quad \text{taken over all open covers } \mathcal{V}.$$

Thus \mathcal{E}_{top} counts, in some sense, the growth rate in the number of T -orbits of length n .

Evidently, topological entropy is an isomorphism invariant. Two continuous maps $T:X \rightarrow X$ and $S:Y \rightarrow Y$ are **topologically conjugate** (as *isomorphism* is called in this category) if there exists a homeomorphism $\psi:X \rightarrow Y$ with $\psi T = S \psi$.

33: Subadditive Lemma. Consider a sequence $\mathbf{s} = (s_\ell)_1^\infty \subset [-\infty, \infty]$ satisfying $s_{k+\ell} \leq s_k + s_\ell$, for all $k, \ell \in \mathbb{Z}$. Then the following limit exists in $[-\infty, \infty]$, and $\lim_{n \rightarrow \infty} \frac{s_n}{n} = \inf_n \frac{s_n}{n}$. ♦

Topological entropy, or “top-ent” for short, satisfies many of the relations of measure-entropy.

34: Lemma.

a: $\mathcal{V} \leq \mathcal{W}$ implies $\mathcal{H}(\mathcal{V}) \leq \mathcal{H}(\mathcal{W})$ and $\mathcal{E}(T, \mathcal{V}) \leq \mathcal{E}(T, \mathcal{W})$.

b: $\mathcal{H}(\mathcal{V} \vee \mathcal{W}) \leq \mathcal{H}(\mathcal{V}) + \mathcal{H}(\mathcal{W})$.

c: $\mathcal{H}(T(\mathcal{V})) \leq \mathcal{H}(\mathcal{V})$, with equality if T is surjective. Also, $\mathcal{E}(T, \mathcal{V}) \leq \mathcal{H}(\mathcal{V})$,

^{♥17}Because we only work on a compact space, we can omit “finite”. Some generalizations of topological entropy to non-compact spaces require that only *finite* open-covers be used; see [37].

d: In (31), the $\lim_{n \rightarrow \infty} \frac{1}{n} \cdot \mathcal{H}(\mathcal{V}_{[0..n]})$ exists.

e: Suppose T is a homeomorphism. Then $\mathcal{E}(T^{-1}, \mathcal{V}) = \mathcal{E}(T, \mathcal{V})$, for each cover \mathcal{V} . Consequently, $\mathcal{E}_{\text{top}}(T^{-1}) = \mathcal{E}_{\text{top}}(T)$.

f: Suppose \mathfrak{C} is a collection of covers such that: For each cover \mathcal{W} , there exists a $\mathcal{V} \in \mathfrak{C}$ with $\mathcal{V} \succcurlyeq \mathcal{W}$. Then $\mathcal{E}_{\text{top}}(T)$ equals the supremum of $\mathcal{E}_{\text{top}}(T, \mathcal{V})$, just taken over those $\mathcal{V} \in \mathfrak{C}$.

g: For all $\ell \in \mathbb{N}$: $\mathcal{E}_{\text{top}}(T^\ell) = \ell \cdot \mathcal{E}_{\text{top}}(T)$. ♦

Proof of (c). Let $\mathcal{C} \leq \mathcal{V}$ be a min-cardinality subcover. Then $T\mathcal{C}$ is a subcover of $T\mathcal{V}$. So $\text{Card}(T\mathcal{V}) \leq |T\mathcal{C}| = |\mathcal{C}|$.

As for entropy, inequality (b) and the foregoing give $\mathcal{H}(\mathcal{V}_{[0..n]}) \leq \mathcal{H}(\mathcal{V}) \cdot n$. ♦

Proof of (d). Set $s_n := \mathcal{H}(\mathcal{V}_{[0..n]})$. Then $s_{k+\ell} \leq s_k + \mathcal{H}(T^\ell(\mathcal{V}_{[0..k]})) \leq s_k + s_\ell$, by (b) and (c), and so the Subadditive Lemma (33), applies. ♦

Proof of (g). WLOG, $\ell = 3$. Given \mathcal{V} a cover, triple it to $\widehat{\mathcal{V}} := \mathcal{V} \cap T\mathcal{V} \cap T^2\mathcal{V}$; so $\bigvee_{j \in [0..N]} [T^3]^j(\widehat{\mathcal{V}}) = \bigvee_{i \in [0..3N]} T^i(\mathcal{V})$. Thus $\mathcal{H}(T^3, \widehat{\mathcal{V}}, N) = \mathcal{H}(T, \mathcal{V}, 3N)$, extending notation. Part (d) and sending $N \rightarrow \infty$, gives $\mathcal{E}(T^3, \widehat{\mathcal{V}}) = 3 \cdot \mathcal{H}(T, \mathcal{V})$.

Lastly, take covers such that $\mathcal{E}(T^3, \mathcal{C}^{(k)}) \rightarrow \mathcal{E}_{\text{top}}(T^3)$ and $\mathcal{E}(T, \mathcal{D}^{(k)}) \rightarrow \mathcal{E}_{\text{top}}(T)$, as $k \rightarrow \infty$. Define $\mathcal{V}^{(k)} := \mathcal{C}^{(k)} \vee \mathcal{D}^{(k)}$. Apply the above to $\mathcal{V}^{(k)}$, then send $k \rightarrow \infty$. ♦

Using a metric

From now on, *our space is a compact metric space* (X, \mathfrak{d}) . Dinaburg [33], and Bowen [34],[35], gave alternative, equivalent, definitions of **top-ent**, in the compact metric-space case, that are often easier to work with than covers. Bowen gave a definition also when X is not compact, ^{♥18} see [35] and [22, chap. 7].

Metric preliminaries. An ε -**ball-cover** comprises finitely many balls, all of radius ε . Since our space is compact, every cover \mathcal{V} has a **Lebesgue number** $\varepsilon > 0$. I.e., for each $z \in X$, the $\text{Bal}(z, \varepsilon)$ lies entirely inside at least one \mathcal{V} -patch. (In particular, there is an ε -ball-cover which *refines* \mathcal{V} .) Let $\text{LEB}(\mathcal{V})$ be the supremum of the Lebesgue numbers. Courtesy (34f) we can

35: Fix a “universal” list $\mathcal{V}^{(1)} \leq \mathcal{V}^{(2)} \leq \dots$, with $\mathcal{V}^{(k)}$ a $\frac{1}{k}$ -ball-cover. For every $T: X \rightarrow X$, then, the $\lim_k \mathcal{E}(T, \mathcal{V}^{(k)})$ computes $\mathcal{E}_{\text{top}}(T)$.

^{♥18}When X is not compact, the definitions need not coincide; e.g [37]. And topologically-equivalent metrics, but which are not uniformly equivalent, may give the same T different entropies, [22, P. 171].

An ε -microscope. Three notions are useful in examining a metric space (X, m) at scale ε . Subset $A \subset X$ is an **ε -separated-set**, if $m(z, z') \geq \varepsilon$ for all distinct $z, z' \in A$. Subset $F \subset X$ is **ε -spanning** if $\forall x \in X, \exists z \in F$ with $m(x, z) < \varepsilon$.

Lastly, a cover \mathcal{V} is **ε -small** if $\text{Diam}(U) < \varepsilon$, for each $U \in \mathcal{V}$.

You take the High Road and I'll take the Low Road. There are several routes to computing $\mathcal{E}_{\text{top}}(T)$, some via maximization, others, minimization. Our foregoing discussion computed $\mathcal{E}_{\text{top}}(T)$ by a family of *sizes* $f_k(n) = f_k^T(n)$, depending on a parameter k which specifies the fineness of scale. (In (35), this k is an integer; in the original definition, an open cover.) Define two numbers:

$$36: \quad \widehat{\mathcal{L}}^f(k) := \limsup_{n \rightarrow \infty} \quad \text{and} \quad \underline{\mathcal{L}}^f(k) := \liminf_{n \rightarrow \infty} \quad \text{of} \quad \frac{1}{n} \log(f_k(n)).$$

Finally, let $\mathcal{E}^f(T) := \sup_k \widehat{\mathcal{L}}^f(k)$. If the limit exists in (36) then we write $\mathcal{L}^f(k)$ for the common value.

The A-K-M definition used the size $f_{\mathcal{V}}(n) := \text{Card}(\mathcal{V}_{[0..n)})$, where

$$\text{Card}(\mathcal{W}) := \text{Minimum cardinality of a subcover from } \mathcal{W}.$$

Here are three metric-space sizes $f_{\varepsilon}(n)$:

$$\text{Sep}(n, \varepsilon) := \text{Maximum cardinality of a } d_n\text{-}\varepsilon\text{-separated set.}$$

$$\text{Spn}(n, \varepsilon) := \text{Minimum cardinality of a } d_n\text{-}\varepsilon\text{-spanning set.}$$

$$\text{Cov}(n, \varepsilon) := \text{Minimum cardinality of a } d_n\text{-}\varepsilon\text{-small cover.}$$

These use a list $(d_n)_{n=1}^{\infty}$ of progressively finer metrics on X , where

$$d_N(x, y) := \max_{j \in [0..N)} d(T^j(x), T^j(y)).$$

37: All-Roads-lead-to-Rome Theorem. Fix ε and let \mathcal{W} be any d - ε -small cover. Then

$$i: \forall n: \text{Cov}(n, 2\varepsilon) \leq \text{Spn}(n, \varepsilon) \leq \text{Sep}(n, \varepsilon) \leq \text{Card}(\mathcal{W}_{[0..n)}).$$

$$ii: \text{Take a cover } \mathcal{V} \text{ and a } \delta < \text{LEB}(\mathcal{V}). \text{ Then } \forall n: \text{Card}(\mathcal{V}_{[0..n)}) \leq \text{Cov}(n, \delta).$$

$$iii: \text{The limit } \mathcal{L}^{\text{Cov}}(\varepsilon) = \lim_n \frac{1}{n} \log(\text{Cov}(n, \varepsilon)) \text{ exists in } [0, \infty).$$

$$iv: \mathcal{E}^{\text{Sep}}(T) = \mathcal{E}^{\text{Spn}}(T) = \mathcal{E}^{\text{Cov}}(T) = \mathcal{E}^{\text{Card}}(T) \stackrel{\text{by defn}}{=} \mathcal{E}_{\text{top}}(T). \quad \diamond$$

Pf of (i). Take $F \subset X$, a min-cardinality d_n - ε -spanning set. So $\bigcup_{z \in F} D_z = X$, where

$$D_z := d_n\text{-Bal}(z, \varepsilon) \stackrel{\text{note}}{=} \bigcap_{j=0}^{n-1} T^{-j}(\text{Bal}(T^j z, \varepsilon)).$$

This $\mathcal{D} := \{D_z\}_z$ is a cover, and it is d_n - 2ε -small. Thus $\text{Cov}(n, 2\varepsilon) \leq |\mathcal{D}| = |F|$.

For any metric, a *maximal* ε -separated-set is automatically ε -spanning; adjoin a putative unspanned point to get a larger separated set.

Let A be a max-cardinality d_n - ε -separated set. Take \mathcal{C} , a min-cardinality subcover of $\mathcal{W}_{[0..n]}$. For each $z \in A$, pick a \mathcal{C} -patch $C_z \ni z$. Could some pair $x, y \in A$ pick the same C ? Well, write $C = \bigcap_{j=0}^{n-1} T^{-j}(W_j)$, with each $W_j \in \mathcal{W}$. For every $j \in [0..n)$, then, $d(T^j(x), T^j(y)) \leq \text{Diam}(W_j) < \varepsilon$. Hence $d_n(x, y) < \varepsilon$; so $x = y$. Accordingly, the $z \mapsto C_z$ map is injective, whence $|A| \leq |\mathcal{C}|$. ♦

Pf of (ii). Choose a min-cardinality d_n - δ -small cover \mathcal{C} . For each $C \in \mathcal{C}$ and $j \in [0..n)$, the d - $\text{Diam}(T^j(C)) < \delta$. So there is a \mathcal{V} -patch $V_{C,j} \supset T^j(C)$. Hence

$$\mathcal{V}_{[0..n)} \stackrel{\text{note}}{\ni} \bigcap_{j=0}^{n-1} T^{-j}(V_{C,j}) \supset C.$$

Thus $\mathcal{V}_{[0..n)} \leq \mathcal{C}$. So $\text{Card}(\mathcal{V}_{[0..n)}) \leq \text{Card}(\mathcal{C}) \leq |\mathcal{C}| = \text{Cov}(n, \delta)$. ♦

Pf of (iii). To upper-bound $\text{Cov}(k+\ell, \varepsilon)$ let \mathcal{V} and \mathcal{W} be min-cardinality ε -small covers, respectively, for metrics d_k and d_ℓ . Then $\mathcal{V} \cap T^\ell(\mathcal{W})$ is a ε -small for $d_{k+\ell}$. Consequently $\text{Cov}(k+\ell, \varepsilon) \leq \text{Cov}(k, \varepsilon) \cdot \text{Cov}(\ell, \varepsilon)$. Thus $n \mapsto \log(\text{Cov}(n, \varepsilon))$ is subadditive. ♦

Pf of (iv). Pick a \mathcal{V} from the list in (35), choose some $2\varepsilon < \text{LEB}(\mathcal{V})$ followed by an ε -small \mathcal{W} from (35). Pushing $n \rightarrow \infty$ gives

$$38: \quad \mathbb{L}^{\text{Card}}(\mathcal{V}) \leq \mathbb{L}^{\text{Cov}}(2\varepsilon) \leq \frac{\widehat{\mathbb{L}}^{\text{Spn}}(\varepsilon) \leq \widehat{\mathbb{L}}^{\text{Sep}}(\varepsilon)}{\underline{\mathbb{L}}^{\text{Spn}}(\varepsilon) \leq \underline{\mathbb{L}}^{\text{Sep}}(\varepsilon)} \leq \mathbb{L}^{\text{Card}}(\mathcal{W}).$$

Now send \mathcal{V} and \mathcal{W} along the (35) list. ♦

Pretension. Topological entropy takes its values in $[0, \infty]$. A useful corollary of (38) can be stated in terms of any $\text{Distance}(\cdot, \cdot)$ which topologizes $[0, \infty]$ as a compact interval.

For each continuous $T: X \rightarrow X$ on a compact metric-space, the $\text{Distance}(\widehat{\mathbb{L}}^{\text{Sep}}(\varepsilon), \underline{\mathbb{L}}^{\text{Sep}}(\varepsilon))$ goes to zero as $\varepsilon \searrow 0$. Consequently, we can pretend that the

$$39: \quad \mathbb{L}^{\text{Sep}}(\varepsilon) = \lim_{n \rightarrow \infty} \frac{1}{n} \log(\text{Sep}(n, \varepsilon))$$

limit exists, in arguments that subsequently send $\varepsilon \searrow 0$. Ditto for $\mathbb{L}^{\text{Spn}}(\varepsilon)$.

We'll use this during the proof of the Variational Principle. But first, here are two entropy computations which illustrate the efficacy in having several characterizations of top-ent.

E.g: $\mathcal{E}_{\text{top}}(\text{Isometry}) = 0$. Suppose $(T : X, d)$ is a distance-preserving map of a compact metric-space. Fixing ε , a set is d_n - ε -separated IFF it is d - ε -separated. Thus $\text{Sep}(n, \varepsilon)$ does *not* grow with n . So each $\widehat{L}^{\text{Sep}}(\varepsilon)$ is zero. \square

E.g: Topological Markov Shifts. Imagine ourselves back in the days when computer data is stored on large reels of fast-moving magnetic tape. One strategy to maximize the density of binary data stored is to *not* put timing-marks (which take up space) on the tape. This has the defect that if we write, say, 577 consecutive 1-bits, the tape-reader may erroneously count 578 copies of 1. We sidestep this flaw by first encoding our data so as to avoid the 11⁵⁷⁷ 1 word, then writing to tape.

Generalize this to a finite alphabet Q and a finite list \mathcal{F} of disallowed Q -words. Extend each word to a common length $K+1$; now $\mathcal{F} \subset Q^{K+1}$. The resulting “ K -step TMS” (*topological Markov shift*) is the shift on the set of doubly- ∞ Q -names having no substring in \mathcal{F} . In the above magnetic-tape example, $K = 576$. Making it more realistic, suppose some string of zeros, say $00^{574}0$, is also forbidden^{♥19} Extending to length 577, we get $2^3=8$ new disallowed words of form $00^{574}0b_1b_2b_3$.

We *recode* to a 1-step TMS (just called a TMS or a *subshift of finite type*) over the alphabet $P := Q^K$. Each outlawed Q -word $w_0w_1 \cdots w_K$ engenders a length-2 forbidden P -word $(w_0, \dots, w_{K-1})(w_1, \dots, w_K)$. The resulting TMS is topologically conjugate to the original K -step. The *allowed* length-2 words can be viewed as the edges in a directed-graph and the set of points $x \in X$ is the set of doubly- ∞ paths through the graph. Once trivialities removed, this X is a Cantor set and the shift $T:X \rightarrow X$ is a homeomorphism.

The Golden Shift. As the simplest example, suppose our magnetic-tape is constrained by the Markov graph, FIG. 40, that we studied measure-theoretically in (26).

We want to store the text of *The Declaration of Independence* on our magnetic tape.^{♥20} Imagining that English is a stationary process, we’d like to encode English into this Golden TMS as efficiently as possible. We seek a shift-invariant measure μ on X_{Gold} of *maximum entropy*, should such exist.

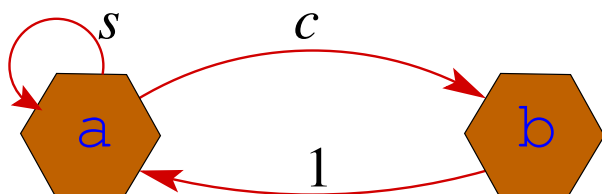


FIG. 40: Ignoring the labels on the edges, for the moment, the **Golden shift**, T , acts on the space of doubly-infinite paths through this graph. The space can be represented as a subset $X_{\text{Gold}} \subset \{a, b\}^{\mathbb{Z}}$, namely, the set of sequences with no two consecutive b letters.

View $P=\{a, b\}$ as the time-zero partition on X_{Gold} ; that is, name $x = \dots x_{-1}x_0x_1x_2 \dots$, is in

^{♥19}Perhaps the 0-bad-length, 574, is shorter than the 1-bad-length because, say, 0s take less tape-space than 1s and so –being written more densely– cause ambiguity sooner.

^{♥20}... which, by Rights, *should* be stored as a Bernoulli process...

atom b IFF letter x_0 is “ b ”. Any measure μ gives conditional probabilities

$$\begin{aligned}\mu(a | a) &=: s, & \mu(b | a) &=: c, \\ \mu(a | b) &\stackrel{\text{note}}{=} 1, & \mu(b | b) &\stackrel{\text{note}}{=} 0.\end{aligned}$$

But recall, $\mathcal{E}(T) = \mathcal{H}(P_1 | P_{[-\infty..0)}) \leq \mathcal{H}(P_1 | P_0)$. So among all measures that make the conditional distribution $P|a$ equal (s, c) , the *unique* one maximizing entropy is the (s, c) -Markov-process. Its entropy, derived in (27), is

$$41: \quad f(s) := \frac{1}{2-s} \cdot \mathcal{H}(s, 1-s) = \frac{-1}{2-s} \cdot [s \log(s) + [1-s] \log(1-s)].$$

Certainly $f(0) = f(1) = 0$, so f ’s maximum occurs at the (it turns out) *unique* point \widehat{s} where the derivative $f'(\widehat{s})$ equals zero. This $\widehat{s} = \frac{-1+\sqrt{5}}{2}$. Plugging in, the maximum entropy supportable by the Golden Shift is

$$42: \quad \text{MaxEnt} = \frac{2}{5-\sqrt{5}} \cdot \left[\frac{-1+\sqrt{5}}{2} \log\left(\frac{2}{-1+\sqrt{5}}\right) + \frac{3-\sqrt{5}}{2} \log\left(\frac{2}{3-\sqrt{5}}\right) \right].$$

Exponentiating, the number of μ -typical n -names grows like G^n , where

$$42': \quad G = \left[\frac{2}{-1+\sqrt{5}} \right]^{\frac{-1+\sqrt{5}}{5-\sqrt{5}}} \cdot \left[\frac{2}{3-\sqrt{5}} \right]^{\frac{3-\sqrt{5}}{5-\sqrt{5}}}.$$

This expression^{♥21} looks unpleasant to simplify –it isn’t even obviously an algebraic number– and yet topological entropy will reveal its familiar nature. This, because the **Variational Principle** (proved in the next section) says that the **top-ent** of a system is the supremum of measure-entropies supportable by the system.

Top-ent of the Golden Shift. For a moment, let’s work more generally on an arbitrary subshift (a closed, shift-invariant subset) $X \subset \mathbb{Q}^{\mathbb{Z}}$, where \mathbb{Q} is a finite alphabet. Here, the transformation is always the shift –but the *space* is varying– so agree to refer to the **top-ent** as $\mathcal{E}_{\text{top}}(X)$. Let $\text{Names}_X(n)$ be the number of distinct words in the set $\{x \upharpoonright_{[0..n)} \mid x \in X\}$. Note that a metric inducing the product-topology on $\mathbb{Q}^{\mathbb{Z}}$ is

$$43: \quad d(x, x') := \frac{1}{1+|m|}, \text{ for the smallest } |m| \text{ with } x_m \neq x'_m.$$

44: Lemma. *Consider a subshift X . Then the $\lim_{n \rightarrow \infty} \frac{1}{n} \log(\text{Names}_X(n))$ exists in $[0, \infty]$, and equals $\mathcal{E}_{\text{top}}(X)$.* ♦

Proof. With $\varepsilon \in (0, 1)$ fixed, two n -names are d_n - ε -separated IFF they are not the same name. Hence $\text{Sep}(n, \varepsilon) = \text{Names}_X(n)$. ♦

^{♥21} A popular computer-algebra-system was not, at least under my inexpert tutelage, able to simplify this. However, once top-ent gave the correct answer, it was able to detect the *equality*.

To compute $\mathcal{E}_{\text{top}}(X_{\text{Gold}})$, declare that a word is “golden” if it appears in some $x \in X_{\text{Gold}}$. Each $[n+1]$ -golden word ending in **a** has form wa , where w is n -golden. An $[n+1]$ -golden word ending in **b**, must end in **ab** and so has form wab , where w is $[n-1]$ -golden. Summing up,

$$\text{Names}_{X_{\text{Gold}}}(n+1) = \text{Names}_{X_{\text{Gold}}}(n) + \text{Names}_{X_{\text{Gold}}}(n-1).$$

This is the Fibonacci recurrence, and indeed, these are the Fibonacci numbers, since $\text{Names}_{X_{\text{Gold}}}(0) = 1$ and $\text{Names}_{X_{\text{Gold}}}(1) = 2$. Consequently, we have that

$$\text{Names}_{X_{\text{Gold}}}(n) \sim \text{Const} \cdot \lambda^n, \quad \text{where } \lambda = \frac{1+\sqrt{5}}{2} \text{ is the Golden Ratio.}$$

So the sesquipedalian number G from (42') is simply λ , and $\mathcal{E}_{\text{top}}(X_{\text{Gold}}) = \log(\lambda)$.

Since $\log(\lambda) \approx 0.694$, each thousand bits written on tape (subject to the “no **bb** substrings” constraint) can carry at most 694 bits of information.

Top-ent of a general TMS. A (finite) digraph G engenders a TMS $T: X_G \rightarrow X_G$, as well as a $\{0, 1\}$ -valued adjacency matrix $A=A_G$, where $a_{i,j}$ is the number of directed-edges from state i to j . (Here, each $a_{i,j}$ is 0 or 1.) The (i, j) -entry in power A^n is automatically the number of length- n paths from i to j . Employing the matrix-norm $\|M\| := \sum_{i,j} |m_{i,j}|$, then,

$$\|A^n\| = \text{Names}_X(n).$$

Happily Gelfand's formula^{♥22} applies: For an arbitrary (square) complex matrix,

$$45: \quad \lim_{n \rightarrow \infty} \|A^n\|^{1/n} = \text{SpecRad}(A).$$

This righthand side, the **spectral radius** of A , means the maximum of the absolute values of A 's eigenvalues. So the top-ent of a TMS is thus the

$$46: \quad \mathcal{E}_{\text{top}}(X_G) = \text{SpecRad}(A_G) := \text{Max}\{|e| \mid e \text{ is an eigenvalue of } A_G\}.$$

The (a, b) -adjacency matrix of Fig. 40 is $\begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$, whose eigenvalues are λ and $-1/\lambda$.

Labeling edges. Interpret $(s, c, 1)$ simply as edge-labels in (40). The set of doubly- ∞ paths can also be viewed as a subset $Y_{\text{Gold}} \subset \{s, c, 1\}^{\mathbb{Z}}$, and it too is a TMS. The shift on Y_{Gold} is conjugate (topologically isomorphic) to the shift on X_{Gold} , so they *a fortiori* have the same top-ent, $\log(\lambda)$. The $(s, c, 1)$ -adjacency matrix is $\begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$. Its $|\cdot|$ -largest eigenvalue is still λ , as it must.

Now we make a new graph. We modify (40) by manufacturing a total of two s -edges, seven c -edges, and three edges $1_1, 1_2, 1_3$. Give these $2+7+3$ edges twelve distinct labels. We *could* compute the resulting TMS-entropy from the corresponding 12×12 adjacency matrix. Alternatively, look at the (a, b) -adjacency matrix $A := \begin{bmatrix} 2 & 7 \\ 3 & 0 \end{bmatrix}$. The roots of its characteristic polynomial are $1 \pm \sqrt{22}$. Hence \mathcal{E}_{top} of this 12-symbol TMS is $\log(1 + \sqrt{22})$. □

^{♥22}See [54, 10.13] or [53, Spectral radius].

The variational principle

Let $\mathfrak{M} := \mathfrak{M}(X, \mathbf{d})$ be the set of Borel probability measures, and $\mathfrak{M}(T) := \mathfrak{M}(T : X, \mathbf{d})$ the set of T -invariant $\mu \in \mathfrak{M}$. Assign

$$\text{EntSup}(T) := \sup \{ \mathcal{E}_\mu(T) \mid \mu \in \mathfrak{M}(T) \}.$$

47: Variational Principle (Goodson). $\text{EntSup}(T) = \mathcal{E}_{\text{top}}(T)$. ◇

This says that top-ent is the top entropy —*if* there is a measure μ which realizes the supremum. There doesn't have to be. Choose a sequence of metric-systems $(S_k: Y_k, \mathbf{m}_k)$ whose entropies *strictly* increase $\mathcal{E}_{\text{top}}(S_k) \nearrow L$ to some limit in $(0, \infty]$. Let $(S_\infty: Y_\infty, \mathbf{m}_\infty)$ be the identity-map on a 1-point space. Define a new system $(T: X, \mathbf{d})$, where $X := \bigsqcup_{k \in [1 \dots \infty]} Y_k$. Have $T(x) := S_k(x)$, for the unique k with $Y_k \ni x$. As for the metric, on Y_k let \mathbf{d} be a scaled version of \mathbf{m}_k , so that the \mathbf{d} -Diam(Y_k) is less than $1/2^k$. Finally, for points in *distinct* components, $x \in Y_k$ and $z \in Y_\ell$, decree that $\mathbf{d}(x, z) := |2^{-k} - 2^{-\ell}|$. Our T is continuous, and is a homeomorphism if each of the S_k is. Certainly $\mathcal{E}_{\text{top}}(T) = L > \mathcal{E}_{\text{top}}(S_k)$, for every $k \in [1 \dots \infty]$.

If L is *finite* then there is *no* measure μ of maximal entropy; for μ must give mass to some Y_k ; this pulls the entropy below L , since there are no compensatory components with entropy exceeding L .

In contrast, when $L = \infty$ then there *is* a maximal-entropy measure (put mass $1/2^j$ on some component Y_{k_j} , where $k_j \nearrow \infty$ swiftly); indeed, there are continuum-many maximal-entropy measures. But there is no^{♥23} *ergodic* measure of maximal entropy.

For a concrete $L=\infty$ example, let S_k be the shift on $[1 \dots k]^{\mathbb{Z}}$.

Topology on \mathfrak{M} . Let's arrange our tools for establishing the Variational Principle. I follow Misiurewicz's proof, adapted from the presentations in [22] and [11].

Equip \mathfrak{M} with the *weak-** topology.^{♥24} An $A \subset X$ is μ -*nice* if its topological boundary $\partial(A)$ is μ -null. And a *partition* is μ -*nice* if each atom is.

48: Prop'n. *If $\alpha_L \rightarrow \mu$ and $A \subset X$ is μ -nice, then $\alpha_L(A) \rightarrow \mu(A)$.* ◇

Proof. Define operator $\mathcal{U}(D) := \limsup_L \alpha_L(D)$. It suffices to show that $\boxed{\mathcal{U}(A) \leq \mu(A)}$. For since A^c is μ -nice too, then $\mathcal{U}(A^c) \leq \mu(A^c)$. Thus $\lim_L \alpha_L(A)$ exists, and equals $\mu(A)$.

Because $C := \overline{A}$ is closed, the continuous functions $f_N \searrow \mathbf{1}_C$ pointwise, where

$$f_N(x) := 1 - \text{Min}(N \cdot \mathbf{d}(x, C), 1).$$

^{♥23}The ergodic measures are the extremepoints of $\mathfrak{M}(T)$; call them $\mathfrak{M}_{\text{Erg}}(T)$. This $\mathfrak{M}(T)$ is the set of barycenters obtained from Borel probability measures on $\mathfrak{M}_{\text{Erg}}(T)$; see [53, Krein-Milman_theorem, Choquet_theory]. In this instance, what explains the failure to have an *ergodic* maximal-entropy measure? Let μ_k be an invariant ergodic measure on Y_k . These measures *do* converge to the one-point (ergodic) probability measure, μ_∞ , on Y_∞ . But the map $\mu \mapsto \mathcal{E}_\mu(T)$ is not continuous at μ_∞ .

^{♥24}Measures $\alpha_L \rightarrow \mu$ IFF $\int f d\alpha_L \rightarrow \int f d\mu$, for each continuous $f: X \rightarrow \mathbb{R}$. This metrizable topology makes \mathfrak{M} compact. Always, $\mathfrak{M}(T)$ is a non-void compact subset; see [Petersen, §6].

By the Monotone Convergence theorem, then, $\int f_N d\mu \xrightarrow{N} \mu(C)$. And $\mu(C) = \mu(A)$, since A is nice. Fixing N , then, it suffices to establish $\mathcal{U}(A) \leq \int f_N d\mu$. But f_N is continuous, so

$$\int f_N d\mu = \limsup_{L \rightarrow \infty} \int f_N d\alpha_L \geq \limsup_{L \rightarrow \infty} \int \mathbf{1}_A d\alpha_L = \mathcal{U}(A). \quad \blacklozenge$$

49: Corollary. Suppose $\alpha_L \rightarrow \mu$ and partition \mathbf{P} is μ -nice. Then $\mathcal{H}_{\alpha_L}(\mathbf{P}) \rightarrow \mathcal{H}_\mu(\mathbf{P})$. \blacklozenge

The **diameter** of partition \mathbf{P} is $\text{Max}_{A \in \mathbf{P}} \text{Diam}(A)$.

50: Prop'n. Take $\mu \in \mathfrak{M}$ and $\varepsilon > 0$. Then there exists a μ -nice partition with $\text{Diam}(\mathbf{P}) < \varepsilon$. \blacklozenge

Proof. Centered at an x , the uncountably many balls $\{\text{Bal}(x, r) \mid r \in (0, \varepsilon)\}$ have disjoint boundaries. So all but countably many are μ -nice; pick one and call it B_x . Compactness gives a finite nice cover, say, $\{B_1, \dots, B_7\}$, at different centers. Then the partition $\mathbf{P} := (A_1, \dots, A_7)$ is nice,^{♥25} where $A_k := B_k \setminus \bigcup_{j=1}^{k-1} B_j$. \blacklozenge

Here is a consequence of Jensen's Inequality.

51: Distropy-averaging Lemma. For $\mu, \nu \in \mathfrak{M}$, a partition \mathbf{R} , and a number $t \in [0, 1]$:

$$t \cdot \mathcal{H}_\mu(\mathbf{R}) + t^c \cdot \mathcal{H}_\nu(\mathbf{R}) \leq \mathcal{H}_{t\mu + t^c\nu}(\mathbf{R}). \quad \blacklozenge$$

Strategy for $\text{EntSup}(T) \geq \mathcal{E}_{\text{top}}(T)$. Choose an $\varepsilon > 0$. For $L = 1, 2, 3, \dots$, take a maximal (L, ε) -separated-set $F_L \subset X$, then define

$$\mathbf{F} = \mathbf{F}_\varepsilon := \limsup_{L \rightarrow \infty} \frac{1}{L} \cdot \log(|F_L|).$$

Let $\varphi_L()$ be the equi-probable measure on F_L ; each point has weight $1/|F_L|$. We will construct our desired invariant measure μ from the Cesàro averages

$$\alpha_L := \frac{1}{L} \cdot \sum_{\ell \in [0..L)} T^\ell \varphi_L,$$

which get more and more invariant.

52: Lemma. Let μ be any weak-* accumulation point of the above $\{\alpha_L\}_1^\infty$. (Automatically, μ is T -invariant.) Then $\mathcal{E}_\mu(T) \geq \mathbf{F}$. Indeed, if \mathbf{Q} is any μ -nice partition with $\text{Diam}(\mathbf{Q}) < \varepsilon$, then $\mathcal{E}_\mu(T, \mathbf{Q}) \geq \mathbf{F}$. \blacklozenge

^{♥25}For any two sets $B, B' \subset X$, the union $\partial B \cup \partial B'$ is a superset of the three boundaries $\partial(B \cup B')$, $\partial(B \cap B')$, $\partial(B \setminus B')$.

Tactics. As usual, $\mathbf{Q}_{[0..N)}$ means $\mathbf{Q}_0 \vee \mathbf{Q}_1 \vee \dots \vee \mathbf{Q}_{N-1}$. Our goal is

$$*: \quad \forall N : \quad \mathbf{F} \stackrel{?}{\leq} \frac{1}{N} \cdot \mathcal{H}_\mu(\mathbf{Q}_{[0..N)}).$$

Fix N and $\mathbf{P} := \mathbf{Q}_{[0..N)}$, and a $\delta > 0$. It suffices to verify: $\forall_{\text{large } L} L \gg N$,

$$52': \quad \frac{1}{L} \log(|F_L|) \stackrel{?}{\leq} \delta + \frac{1}{N} \cdot \mathcal{H}_{\alpha_L}(\mathbf{P}),$$

since this and (49) will prove (*): Pushing $L \rightarrow \infty$ along the sequence that produced μ essentially sends LhS(52') to \mathbf{F} , courtesy (39). And RhS(52') goes to $\delta + \frac{1}{N} \cdot \mathcal{H}_\mu(\mathbf{P})$, by (49), since \mathbf{P} is μ -nice. Descending $\delta \searrow 0$, hands us the needed (*). \square

Proof²⁶ of (52'). Since L is frozen, agree to use φ for the φ_L probability measure.

Our d_L - ε -separated set F_L has at most *one* point in any given atom of $\mathbf{Q}_{[0..L)}$, thereupon

$$\log(|F_L|) = \mathcal{H}_\varphi(\mathbf{Q}_{[0..L)}).$$

Regardless of the “offset” $K \in [0..N)$, we can always fit $C := \lfloor \frac{L-N}{N} \rfloor$ many N -blocks into $[0..L)$. Denote by $\mathcal{G}(K) := [K..K+CN)$, this union of N -blocks, the **good** set of indices. Unsurprisingly, $\mathcal{B}(K) := [0..L) \setminus \mathcal{G}(K)$ is the **bad** index-set. Therefore,

$$53: \quad \mathcal{H}_\varphi(\mathbf{Q}_{[0..L)}) \leq \overbrace{\mathcal{H}_\varphi(\bigvee_{j \in \mathcal{B}(K)} \mathbf{Q}_j)}^{\text{Bad}(K)} + \overbrace{\mathcal{H}_\varphi(\bigvee_{j \in \mathcal{G}(K)} \mathbf{Q}_j)}^{\text{Good}(K)}.$$

Certainly $\text{Bad}(K) \leq 3N \log(|\mathbf{Q}|)$. So $\frac{1}{NL} \sum_{K \in [0..N)} \text{Bad}(K) \leq \frac{3N}{L} \log(|\mathbf{Q}|)$. This is less than δ , since L is large. Applying $\frac{1}{NL} \sum_{K \in [0..N)}$ to (53) now produces

$$54: \quad \frac{1}{L} \cdot \log(|F_L|) \leq \delta + \frac{1}{NL} \sum_K \text{Good}(K).$$

Note $\bigvee_{j \in \mathcal{G}(K)} T^j(\mathbf{Q}) = \bigvee_{c \in [0..C)} T^{K+cN}(\mathbf{P})$. So $\text{Good}(K) \leq \sum_c \mathcal{H}_\varphi(T^{K+cN}(\mathbf{P}))$. This latter, by definition, equals $\sum_c \mathcal{H}_{T^{K+cN}(\varphi)}(\mathbf{P})$. We conclude that

$$\begin{aligned} \frac{1}{NL} \sum_K \text{Good}(K) &\leq \frac{1}{NL} \sum_K \sum_c \mathcal{H}_{T^{K+cN}(\varphi)}(\mathbf{P}) \\ &\leq \frac{1}{NL} \sum_{\ell \in [0..L)} \mathcal{H}_{T^\ell \varphi}(\mathbf{P}), \quad \text{by adjoining a few translates of } \mathbf{P}, \\ &\leq \frac{1}{N} \cdot \mathcal{H}_{\alpha_L}(\mathbf{P}), \quad \text{by the Distropy-averaging Lemma, (51),} \end{aligned}$$

since α_L is the average $\frac{1}{L} \sum_\ell T^\ell \varphi$. Thus (54) implies (52'), our goal. \blacklozenge

²⁶The idea is to mostly fill $[0..L)$ with N -blocks, starting with a offset $K \in [0..N)$. Averaging over the offset will create a Cesàro average over each N -block. Averaging over the N -blocks will allow us to compute distropy with respect to the averaged measure, α_L .

Proof of $\text{EntSup}(T) \leq \mathcal{E}_{\text{top}}(T)$. Fix a T -invariant μ . For partition $Q = (B_1, \dots, B_K)$, choose a compact set $A_k \subset B_k$ with $\mu(B_k \setminus A_k)$ small.^{♡27} Letting $D := [\bigsqcup_i A_i]^c$ and $P := (D, A_1, \dots, A_K)$, we can have made $\mathcal{H}(P | Q)$ as small as desired. Courtesy (23b), then, we only need consider partitions of the form that P has.

Open-cover $\mathcal{V} := (U_1, \dots, U_K)$ has patches $U_k := D \cup A_k$. What atoms of, say, $P_{[0..3]}$, can the intersection $U_9 \cap T^{-1}(U_2) \cap T^{-2}(U_5)$ touch? Only the eight atoms

$$(D \text{ or } A_9) \cap T^{-1}(D \text{ or } A_2) \cap T^{-2}(D \text{ or } A_5).$$

Thus $\#P_{[0..n]} \leq 2^n \cdot \#\mathcal{V}_{[0..n]}$. (Here, $\#()$ counts the number of non-void atoms/patches.) So

$$\frac{1}{n} \mathcal{H}_\mu(P_{[0..n]}) \leq 1 + \frac{1}{n} \log(\#\mathcal{V}_{[0..n]}) \leq 1 + 1 + \mathcal{E}_{\text{top}}(T);$$

this last inequality, when n is large. The upshot: $\mathcal{E}_\mu(T) \leq 2 + \mathcal{E}_{\text{top}}(T)$.

Applied to a power T^ℓ , this asserts that $\mathcal{E}_\mu(T^\ell) \leq 2 + \mathcal{E}_{\text{top}}(T^\ell)$. Thus

$$\mathcal{E}_\mu(T) \leq \frac{2}{\ell} + \mathcal{E}_{\text{top}}(T),$$

using (24d) and using (34g). Now coax $\ell \rightarrow \infty$. ♦

Three recent results

Having given an survey of older results in measure-theoretic entropy and in topological entropy, let us end this survey with a brief discussion of a few recent results, chosen from many.

Ornstein-Weiss: Finitely-observable invariant. In a landmark paper [8, 2007], Ornstein and Weiss show that all “finitely observable” properties of ergodic processes are secretly entropy; indeed, they are continuous functions of entropy. This was generalized by Gutman and Hochman [9]; some of the notation below is from their paper.

Here is the setting. Consider an ergodic process, on a non-atomic space, taking on only finitely many values in \mathbb{N} ; let \mathcal{C} be some family of such processes. An **observation scheme** is a metric space (Ω, d) and a sequence of functions $\mathbf{S} = (S_n)_1^\infty$, where S_n maps \mathbb{N}^∞ into Ω . On a point $\vec{x} \in \mathbb{N}^\infty$, the scheme **converges** if

$$55: \quad n \mapsto S_n(x_1, x_2, \dots, x_n)$$

converges in Ω . And on a particular process X , say that \mathbf{S} **converges**, if \mathbf{S} converges on a.e. \vec{x} in X .

A function $J: \mathcal{C} \rightarrow \Omega$ is isomorphism invariant if, whenever the underlying transformations of two processes $X, X' \in \mathcal{C}$ are isomorphic, then $J(X) = J(X')$. Lastly, say that \mathbf{S} “converges to J ”, if for each $X \in \mathcal{C}$, scheme \mathbf{S} converges to the value $J(X)$.

^{♡27}This can be done, since μ is automatically a regular measure.

The work of David Bailey [38, 1976], a student of Ornstein, produced an observation scheme for entropy. The Lempel-Ziv algorithm [43] was another entropy observer, with practical application.

Ornstein and Weiss provided entropy schemes in [41] and [42]. Their recent paper “Entropy is the only finitely-observable invariant” [8, 2007], give a converse, a uniqueness result.

56: Theorem (Ornstein, Weiss). *Suppose J is a finitely observable function, defined on all ergodic finite-valued processes. If J is an isomorphism invariant, then J is a continuous function of the entropy.* \diamond

Yonatan Gutman and Michael Hochman, in [9], significantly extend the Ornstein-Weiss result, by proving that it holds even when the isomorphism invariant, J , is well-defined only on certain subclasses of the set of all ergodic processes. In particular they obtain the following result on three classes of zero-entropy transformations.

57: Theorem (Gutman, Hochman). *Suppose $J()$ is a finitely observable invariant on one of the following classes:*

- i: The Kronecker systems; the class of systems with pure point spectrum.*
- ii: The zero-entropy mild mixing processes.*
- iii: The zero-entropy strong mixing processes.*

Then $J()$ is constant. \diamond

Entropy of actions of free groups. Consider (G, \mathcal{G}) , a topological group and its Borel field (sigma-algebra). Let $\mathcal{G} \times \mathcal{X}$ be the field on $G \times X$ generated by the two coordinate-subfields. A map

*****: $\psi: G \times X \rightarrow X$ is **measurable** if $\psi^{-1}(\mathcal{X}) \subset \mathcal{G} \times \mathcal{X}$. Use $\psi^g(x)$ for $\psi(g, x)$.

This map $(*)$ is a (measure-preserving) **group action** if $\forall g, h \in G: \psi^g \circ \psi^h = \psi^{gh}$, and each $\psi^g: X \rightarrow X$ is measure preserving.

This encyclopedia article has only discussed entropy for \mathbb{Z} -actions, i.e, when $G = \mathbb{Z}$. The ergodic theorem, our definition of entropy, and large parts of ergodic theory, involve taking averages (of some quantity of interest) over larger and larger “pieces of Time”. In \mathbb{Z} , we typically use the intervals $I_n := [0 .. n)$. When G is $\mathbb{Z} \times \mathbb{Z}$, we might average over squares $I_n \times I_n$.

The *amenable groups* are those which possess, in a certain sense, larger and larger averaging sets. Parts of ergodic theory have been carried over to actions of amenable groups, e.g [49] and [51]. Indeed, much of the Bernoulli theory was extended to certain amenable groups by Ornstein and Weiss, [50].

The stereotypical example of a *non*-amenable group, is a free group (on more than one generator). But recently, Lewis Bowen [10] succeeded in extending the definition of entropy to actions of finite-rank free groups.

58: Theorem (Lewis Bowen). *Let G be a finite-rank free group. Then two Bernoulli G -actions are isomorphic IFF they have the same entropy.* \diamond

The paper introduces a new isomorphism invariant, the “ f invariant”, and shows that, for Bernoulli actions, the f invariant agrees with entropy, that is, with the distropy of the independent generating partition.

Conclusion

Ever since the pioneering work of Shannon, and of Kolmogorov and Sinai, entropy has been front and center as a major tool in Ergodic Theory. Simply *mentioning* all the substantial results in entropy theory would dwarf the length of this encyclopedia article many times over. And, as the above three results (cherry-picked out of many) show, Entropy shows no sign of fading away. . .

§BIBLIOGRAPHY

[1] Citations to other articles in this Encyclopedia are made with the author’s name in double-brackets, e.g., [[Lemanczyk]], [[Petersen]].

Historical

- [2] Rudolf Clausius, *Abhandlungen ueber die mechanische Waermetheorie*, vol. **1**, (F.Vieweg, Braunschweig, 1864); vol. **2**, (1867).
- [3] Claude Elwood Shannon, *A Mathematical Theory of Communication*, Bell Syst. Tech. J., 27, pp. 379–423, pp. 623–656, (1948).
- [4] A.N. Kolmogorov, *A New Metric Invariant of Transitive Automorphisms of Lebesgue Spaces*, Dokl. Akad. Nauk SSSR 119, no. 5, pp. 861–864, (1958).
- [5] Ya. Sinai, *On the Concept of Entropy of a Dynamical System*, Dokl. Akad. Nauk SSSR 124, pp. 768–771, (1959).
- [6] B. McMillan, *The Basic Theorems of Information Theory*, Ann. Math. Stat. 24, pp. 196–219, (1953).
- [7] Roy Adler and Benjamin Weiss, *Entropy, a complete metric invariant for automorphisms of the torus*, Proc. Nat. Acad. Sci. USA, vol. **57**, pp. 1573–1576, (1967).

Recent Results

- [8] D.S. Ornstein and B. Weiss, *Entropy is the only finitely-observable invariant* J. of Modern Dynamics, vol. **1**, (2007), pp. 93–105. <http://www.math.psu.edu/jmd>
- [9] Yonatan Gutman and Michael Hochman, *On processes which cannot be distinguished by finitary observation* to appear in Israel Journal of Mathematics. Preprint at <http://arxiv.org/pdf/math/0608310>

- [10] Lewis Bowen, *A new measure-conjugacy invariant for actions of free groups* Preprint at <http://www.math.hawaii.edu/%7Elpbowen/notes11.pdf>

Ergodic Theory books

- [11] Michael Brin, Garrett Stuck, *Introduction to dynamical systems*, Cambridge University Press, (2002).
- [12] I. Cornfeld, S. Fomin, Ya. Sinai, *Ergodic theory*, Grundlehren der Mathematischen Wissenschaften, 245. Springer-Verlag, New York, (1982).
- [13] Nathaniel A. Friedman, *Introduction to Ergodic Theory*, Van Nostrand Reinhold, (1970).
- [14] Paul R. Halmos, *Lectures on Ergodic Theory*. The Mathematical Society of Japan, (1956).
- [15] Anatole Katok, Boris Hasselblatt, *Introduction to the modern theory of dynamical systems. (With a supplementary chapter by Katok and Leonardo Mendoza)*, Encyclopedia of Mathematics and its Applications, 54. Cambridge University Press, (1995).
- [16] G. Keller, A. Greven and G. Warnecke (eds), *Entropy*, Princeton Series in Applied Mathematics, Princeton University Press (2003).
- [17] Doug Lind, Brian Marcus, *An introduction to symbolic dynamics and coding*, Cambridge University Press, Cambridge, (1995).
- [18] R. Mané. *Ergodic theory and differentiable dynamics*. Ergebnisse der Mathematik und ihrer Grenzgebiete; ser.3, Bd. 8. Springer-Verlag, Berlin, (1987).
- [19] William Parry, *Entropy and Generators in Ergodic Theory*, W.A. Benjamin, (1969).
- [20] Karl Petersen, *Ergodic theory*. Cambridge Univ. Press, Cambridge, (1983).
- [21] Daniel J. Rudolph, *Fundamentals of Measurable Dynamics*, Clarendon Press, (1990).
- [22] Peter Walters, *An introduction to ergodic theory*, Graduate Texts in Mathematics, no. 79, Springer, (1982).
- [23] Ya.G. Sinai. *Topics in ergodic theory*, volume 44 of *Princeton Mathematical Series*. Princeton University Press, Princeton, NJ, (1994).

Differentiable entropy

- [24] Ya.B. Pesin. *Characteristic Lyapunov exponents and smooth ergodic theory*. Russ. Math. Surveys, 32:55–114, (1977).
- [25] L.-S. Young, *Dimension, entropy and Lyapunov exponents*, Ergodic Theory and Dynamical Systems, 2, (1982), no. 1, pp. 109–124.
- [26] F. Ledrappier and L.-S. Young, *The metric entropy of diffeomorphisms*, Ann. of Math., (2) 122, (1985), pp. 509–539 (I) and pp. 540–574 (II).

Finite rank

- [27] S. Ferenczi, *Systems of finite rank*, Colloq. Math., 73, (1997), no. 1, pp. 35–65.

- [28] J.L.F. King, *Joining-rank and the structure of finite rank mixing transformations*, *J. Analyse Math.*, vol. **51**, (1988), pp. 182–227.

Maximal-entropy measures

- [29] J. Buzzi and S. Ruelle, *Large entropy implies existence of a maximal entropy measure for interval maps*, *Discrete Contin. Dyn. Syst.*, 14, (2006), no. 4, pp. 673–688.
- [30] M. Denker, *Measures with maximal entropy*, *Théorie ergodique Actes Journées Ergodiques, Rennes, 1973/1974*, pp. 70–112. *Lecture Notes in Math.*, vol. 532, Springer, Berlin, (1976).
- [31] M. Misiurewicz, *Diffeomorphism without any measure with maximal entropy*, *Bull. Acad. Polon. Sci. Sér. Sci. Math. Astronom. Phys.* 21, (1973), pp. 903–910.

Topological entropy

- [32] R.L. Adler, A.G. Konheim, M.H. McAndrew, *Topological Entropy*, *Transactions of the American Mathematical Society*, vol. **114**, no. 2, (Feb., 1965), pp. 309–319.
- [33] E.I. Dinaburg, *The relation between topological entropy and metric entropy*, *Soviet Math. Dokl.* vol. **11**, (1970), pp. 13–16.
- [34] Rufus Bowen, *Entropy for group endomorphisms and homogeneous spaces*, *Transactions, American Mathematical Society*, vol. **153**, (1971), pp. 401–414, erratum, **181**(1973) pp. 509–510.
- [35] Rufus Bowen, *Topological entropy for noncompact sets*, *Transactions, American Mathematical Society*, vol. **184**, (1973), pp. 125–136.
- [36] Roy Adler, Brian Marcus, *Topological entropy and equivalence of dynamical systems*, *Mem. Amer. Math. Soc.*, vol. **20**, (1979), no. 219.
- [37] Boris Hasselblatt, Zbigniew Nitecki, James Propp, *Topological entropy for non-uniformly continuous maps*, <http://www.citebase.org/abstract?id=oai:arXiv.org:math/0511495>, (2005).

Determinism and Zero-entropy, and entropy observation

- [38] David Bailey. *Sequential schemes for clasifying and predicting ergodic processes*. Stanford University, (1976). Ph.D. Dissertation.
- [39] J.L.F. King, *Dilemma of the sleeping stockbroker*, *The American Mathematical Monthly*, **99**, no. 4, (1992), pp. 335–338.
- [40] S. Kalikow and J.L.F. King, *A countably-valued sleeping stockbroker process*, *Journal of Theoretical Probability*, vol. **7**, no. 4, (1994), pp. 703–708.
- [41] Donald S. Ornstein and Benjamin Weiss, *How sampling reveals a process*. *Ann. Probab.*, 18(3):905–930, (1990).
- [42] Donald S. Ornstein and Benjamin Weiss, *Entropy and data compression schemes*. *IEEE Trans. Inform. Theory*, 39(1):78–83, (1993).
- [43] Jacob Ziv and Abraham Lempel. *A universal algorithm for sequential data compression*. *IEEE Trans. Information Theory*, IT-23(3):337–343, (1977).

Bernoulli Transformations, K-automorphisms, Amenable groups

- [44] L.D. Meshalkin, *A Case of Isomorphism of Bernoulli Schemes*, Dokl. Akad. Nauk SSSR 128, pp. 41–44, (1959).
- [45] Ya.G. Sinai, *A Weak Isomorphism of Transformations Having an Invariant Measure*, Dokl. Akad. Nauk. SSSR 147, pp. 797–800, (1962), *MR 161960*, [Zbl 0205.13501](#).
- [46] Donald S. Ornstein, *Bernoulli Shifts with the Same Entropy Are Isomorphic*, Adv. Math., 5, pp. 337–352, (1970).
- [47] Paul Shields. *The Theory of Bernoulli Shifts*, University of Chicago Press, Chicago and London, (1973).
- [48] Donald S. Ornstein, *Ergodic Theory Randomness and Dynamical Systems*, Yale Math. Monographs, vol. 5, Yale University Press, (1974), *MR 447525*, [Zbl 0296.28016](#).
- [49] D.S. Ornstein and B. Weiss, *The Shannon-Mc Millan-Briman Theorem For A Class Of Amenable Groups*, Isr. J. Math., vol. 44, no. 3, (1983), pp. 53–60.
- [50] D.S. Ornstein and B. Weiss, *Entropy and isomorphism theorems for actions of amenable groups*, *J. Analyse Math.*, vol. 48, (1987) pp. 1–141.

Abramov Formula

- [51] T. Ward and Q. Zhang, *The Abramov-Rohlin entropy addition formula for amenable group actions*, *Monatsh. Math.*, 114, (1992), pp. 317–329.

Miscellaneous

- [52] S.E. Newhouse. *Continuity properties of entropy*, *Annals of Mathematics*, 129:215–235, (1989). Also *Ann. of Math.* 131, (1990), pp. 409–410.
- [53] *Wikipedia*, http://en.wikipedia.org/wiki/pages:Spectral_radius,Information_entropy.
- [54] Walter Rudin, *Functional Analysis*, McGraw-Hill, (1973).

Filename: Problems/Dynamics/Entropy/entr_biblio.latex

As of: Friday 13Jul2007. Typeset: 4May2008 at 0-900:20.