

Lecture 3: Entropy

For each probability preserving map T on a probability space $(\Omega, \mathcal{F}, \mu)$ we will define a number $h_\mu(T)$ called the metric entropy (or Kolmogorov-Sinai entropy).

Then we will interpret it in terms of:

- the size of coarse grained phase space
- the information production rate
- the capacity of T to generate purely random signals

The metric entropy is defined using "dynamics of partitions", so we begin this.

Dynamics of Partitions

Defⁿ. A finite measurable partition of a probability space $(\Omega, \mathcal{F}, \mu)$ is a collection $\alpha = \{A_1, \dots, A_N\}$ s.t.:

- A_1, \dots, A_N are measurable sets (called the "atoms")
- $A_i \cap A_j = \emptyset$ for $i \neq j$
- $A_1 \cup \dots \cup A_N = \Omega$

We let $\alpha(x) :=$ the element of α which contains x

Given two partitions $\alpha = \{A_1, \dots, A_N\}, \beta = \{B_1, \dots, B_M\}$

$$\alpha \vee \beta := \left\{ A_i \cap B_j \mid \begin{array}{l} i = 1, \dots, N \\ j = 1, \dots, M \end{array} \right\}$$

How to think about this object?

① experiment with N possible outcomes

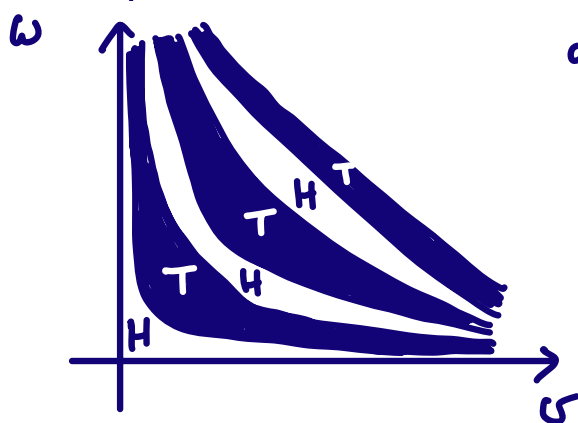
$$A_i = \{\omega \in \Omega : \text{the experiment gives } i^{\text{th}} \text{ outcome}\}$$

$\alpha(\omega) :=$ the outcome when the state is ω

For example, in lecture 1 we saw that when we toss a coin with $\left\{ \begin{array}{l} \text{vertical velocity } v \\ \text{angular velocity } \omega \end{array} \right.$, then

$$\alpha = \left\{ \left\{ (v, \omega) : \text{"heads"} \right\}, \left\{ (v, \omega) : \text{"tails"} \right\} \right\}$$

is the partition



$$\alpha = \left\{ \begin{array}{l} \text{union of} \\ \text{white} \\ \text{bands} \end{array} , \begin{array}{l} \text{union of} \\ \text{blue} \\ \text{bands} \end{array} \right\}$$

Given two "experiments" α, β , if we carry out both at the same time, the possible outcomes are

$$(i, j) = \left(\begin{array}{l} i^{\text{th}} \text{ outcome} \\ \text{in exp. } \alpha \end{array} ; \begin{array}{l} j^{\text{th}} \text{ outcome} \\ \text{in exp. } \beta \end{array} \right)$$

and we obtain $\alpha \vee \beta$:

$$A_i \cap B_j = \left\{ \omega \in \Omega : \begin{array}{l} \alpha \text{ exp. gives } i^{\text{th}} \text{ outcome} \\ \beta \text{ exp. gives } j^{\text{th}} \text{ outcome} \end{array} \right\}$$

② Information: Imagine we don't know ω fully,
 we only know $\alpha(\omega) =$ element of α containing ω .
 With this info

- * can distinguish ω_1, ω_2 in different elements of α
- * cannot ———— in the same elements of α .

For example: Suppose $\Omega = [0,1]^2$

- $\alpha =$ information on 1st binary digit of x
- $\beta =$ ———— 1st ———— y

(0,*)	(1,*)
-------	-------

α

(*,1)
(*,0)

β

(0,1)	(1,1)
(0,0)	(1,0)

$\alpha \vee \beta$

③ Coarse Graining:

$\alpha =$ coarse graining of observers info on $\omega \in \Omega$

We're not coarse graining the system!

We're coarse graining the observer.

— bad for simulations

— good for theory (the dynamics stays the same)

A calculation: Suppose we repeat an "experiment" $\alpha = \{A_1, A_2, \dots\}$ at times $t=0, 1, 2, \dots$.
 What is the information we have at time n ?

time	state	new information	accumulated information
0	ω	$\omega \in A_{i_0} := \alpha(\omega)$	$\omega \in A_{i_0}$
1	$T(\omega)$	$T(\omega) \in A_{i_1} := \alpha(T\omega)$	$\omega \in A_{i_0} \cap T^{-1}A_{i_1}$
2	$T^2(\omega)$	$T^2(\omega) \in A_{i_2} := \alpha(T^2\omega)$	$\omega \in A_{i_0} \cap T^{-1}A_{i_1} \cap T^{-2}A_{i_2}$
\vdots	\vdots	\vdots	\vdots
$n-1$	$T^{n-1}(\omega)$	$T^{n-1}(\omega) \in A_{i_{n-1}} := \alpha(T^{n-1}\omega)$	$\omega \in \bigcap_{j=0}^{n-1} T^{-j}A_{i_j}$

Let $T^{-j}\alpha = \{T^{-j}A_i : A_i \in \alpha\}$, then

$$(T^{-j}\alpha)(\omega) = \alpha(T^j\omega) \quad (\text{check!})$$

The information at time n is coded by the partition

$$\alpha_n = \bigvee_{j=0}^{n-1} T^{-j}\alpha = \left\{ \bigcap_{j=0}^{n-1} T^{-j}A_{i_j} : A_{i_0}, \dots, A_{i_{n-1}} \in \alpha \right\}.$$

Notice that the information grows
 the cardinality of the partition grows
 but the elements ——— " ——— decrease

Example: $T: [0,1] \rightarrow [0,1]$, $T(\omega) = 2\omega \pmod{1}$.

In binary expansions, $2 \times 0.\omega_1\omega_2\omega_3\cdots = \omega_1.\omega_2\omega_3\cdots$ so

$$T(0.\omega_1\omega_2\omega_3\cdots) = 0.\omega_2\omega_3\cdots.$$

Suppose we can only measure the first binary digit:

$$\alpha = \left\{ [0, 1/2); [1/2, 1] \right\}$$

time	state	new info	total info
0	$\omega = 0.\omega_1\omega_2\cdots$	ω_1	ω_1
1	$T(\omega) = 0.\omega_2\omega_3\cdots$	ω_2	ω_1, ω_2
\vdots	\vdots	\vdots	\vdots
$n-1$	$T^{n-1}(\omega) = 0.\omega_n\omega_{n+1}\cdots$	ω_{n-1}	$\omega_1, \dots, \omega_{n-1}$

Thus: $\alpha_n = \bigcup_{i=0}^{n-1} T^{-i} \alpha =$ partition into 2^n dyadic intervals of length $1/2^n$

The Entropy of a Partition

Defⁿ. Let $\alpha = \{A_1, \dots, A_N\}$ be a finite measurable partition of a probability space $(\Omega, \mathcal{F}, \mu)$.

$$H_\mu(\alpha) := - \sum_{A \in \alpha} \mu(A) \ln \mu(A)$$

$$h_\mu(T, \alpha) := \lim_{n \rightarrow \infty} \frac{1}{n} H_\mu(\alpha_n), \quad \alpha_n = \bigvee_{j=0}^{n-1} T^{-j} \alpha.$$

$$h_\mu(T) := \sup \left\{ h_\mu(T, \alpha) : \alpha \text{ a finite measurable partition of } \Omega \right\}$$

Thm (Krieger): If $h_\mu(T) < \infty$, then there exist finite partitions α s.t. $h_\mu(T, \alpha) = h_\mu(T) = \text{maximal possible.}$

Example: Let

- $T(\omega) = 2\omega \text{ mod } 1$ on $[0, 1]$
- $\mu = \text{length measure}$ (Lebesgue measure)
- $\alpha = \left\{ [0, \frac{1}{2}), [\frac{1}{2}, 1) \right\}$. Then:

$\alpha_n = n^{\text{th}}$ dyadic partition

$$H_\mu(\alpha_n) = \underbrace{2^n}_{\# \text{ intervals}} \times \underbrace{\left(-\frac{1}{2^n} \ln \frac{1}{2^n} \right)}_{2^{-n} \cdot n \ln 2} = n \ln 2.$$

$$\text{So } h_\mu(T, \alpha) = \lim_{n \rightarrow \infty} \frac{1}{n} H_\mu(\alpha_n) = \ln 2.$$

It can be shown that $h_\mu(T, \alpha) = h_\mu(T)$.

Entropy and "the Volume of Coarse Grained Space"

Think of $\alpha_n = \{\alpha_n(\omega) : \omega \in \Omega\}$ as of a coarse-graining of the observer's information on Ω .

Shannon-McMillan-Breiman Thm: Let T be an ergodic probability preserving map on a probability space $(\Omega, \mathcal{F}, \mu)$. For every finite measurable partition α ,

(1) for μ -a.e. $\omega \in \Omega$,

$$-\frac{1}{n} \log \mu(\alpha_n(\omega)) \xrightarrow{n \rightarrow \infty} h_\mu(T, \alpha)$$

(2) for every $\epsilon > 0$, for all n large enough we can decompose

$$\Omega = \Omega_{\text{main}} \cup \Omega_{\text{negligible}}$$

so that:

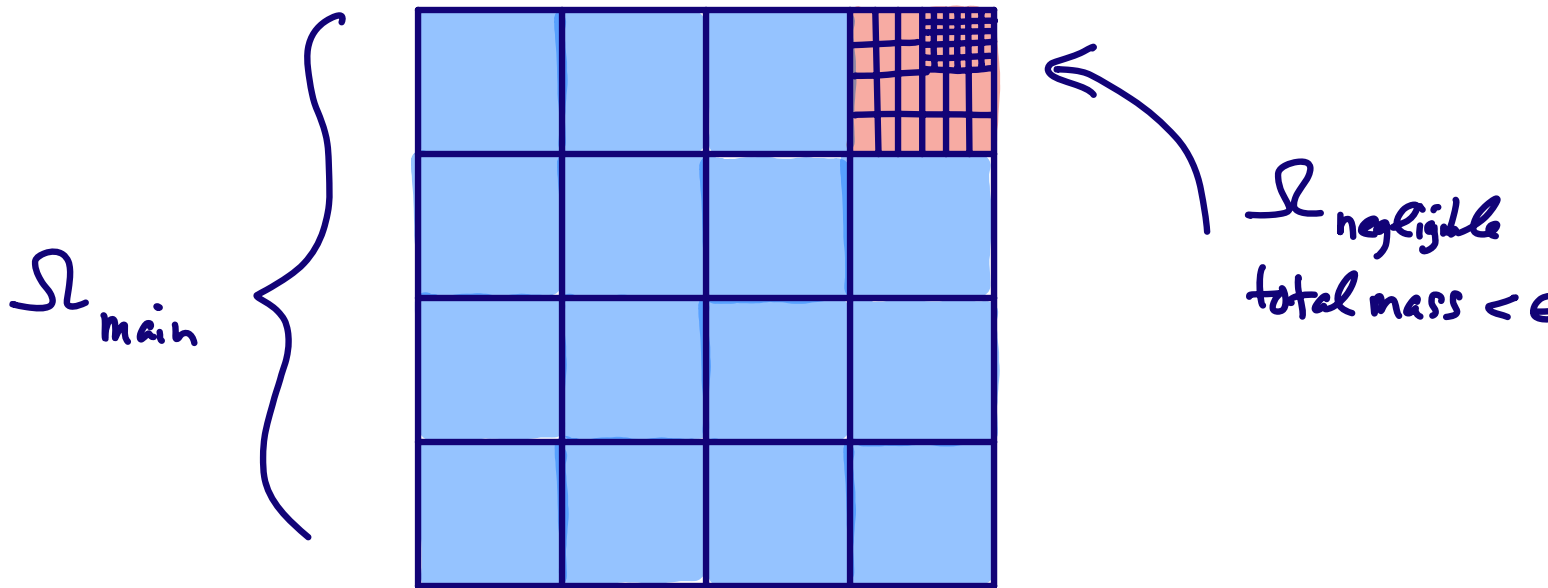
(a) $\mu(\Omega_{\text{negligible}}) < \epsilon$

(b) $\Omega_{\text{main}} = \text{union of } e^{n(h_\mu(T, \alpha) \pm \epsilon)}$

elements of α_n , each of size $e^{-n(h_\mu(T, \alpha) \pm \epsilon)}$

Thus the cardinality of coarse grained space is approximately $\exp(n \cdot (h_f(\tau, \epsilon) \pm \epsilon))$.

"Approximately" = (after removing a set of measure $< \epsilon$)



Entropy and Information Production Rate

We wish to quantify the information content $I_\mu(E)$ of the statement " ω belongs to E ."

Axioms:

(a) $I_\mu(E) =$ decreasing function of $\mu(E)$

(b) $I_\mu(E) =$ continuous function of $\mu(E)$

(c) If A, B are independent, i.e. $\mu(A \cap B) = \mu(A)\mu(B)$, then $I_\mu(A \cap B) = I_\mu(A) + I_\mu(B)$.

All such functions are positively proportional to

$$I_{\mu}(E) := -\ln \mu(E)$$

(in computer science, people use \log_2).

Notice that

$$\begin{aligned} H_{\mu}(\alpha) &= -\sum_{A \in \alpha} \mu(A) \ln \mu(A) \\ &= \int_{\Omega} \sum_{A \in \alpha} \mathbb{1}_A(\omega) \cdot I_{\mu}(A) \, d\mu \\ &= \int_{\Omega} I_{\mu}(\alpha)(\omega) \, d\mu \end{aligned}$$

where $I_{\mu}(\alpha)(\omega) = I_{\mu}(\alpha(\omega))$. Observe:

the information we gain in performing

- $I_{\mu}(\alpha)(\omega) =$ the experiment when the system is at state ω

- $H_{\mu}(\alpha) =$ mean information we gain when performing the experiment (over all $\omega \in \Omega$)

If we repeat an experiment α at times $t=0, 1, \dots$ then at time $n-1$ we know $\alpha_n(\omega)$, and the amount of information we have is $I_\mu(\alpha_n(\omega)) = -\log \mu(\alpha_n(\omega))$.

Corollary. Under the conditions of the Shannon-McMillan-Breiman theorem, for μ -a.e. ω

$$\lim_{n \rightarrow \infty} \frac{1}{n} I_\mu(\alpha_n(\omega)) = h_\mu(T, \alpha)$$

"information production rate"

Entropy and Unpredictability

Define $H_\mu(\alpha | \beta) := H_\mu(\alpha \vee \beta) - H_\mu(\beta)$
 = mean additional information in α given we already know β .

Rokhlin Formula: Under the assumptions of SMS thm, for every finite measurable partition α ,

$$h_\mu(T, \alpha) = \lim_{n \rightarrow \infty} H_\mu(T^{-n} \alpha | \alpha_n)$$

learning $\alpha(T^n \omega)$
given that we already know $\alpha(\omega), \alpha(T\omega), \dots, \alpha(T^{n-1}\omega)$.

If $h_\mu(T, \alpha) > 0$, we still have much to learn, even at time $n \gg 1$!

Entropy and Deterministic Chaos

Deterministic Chaos: The ability of deterministic systems to produce behavior which appears random.

The "most random stochastic process":

Defⁿ. Let (p_1, \dots, p_N) be a probability vector, i.e. $0 \leq p_i \leq 1$ and $p_1 + \dots + p_N = 1$.

A Bernoulli process $\mathcal{B}(p_1, \dots, p_N)$ is a bi-infinite sequence of independent random variables $\{X_i\}_{i \in \mathbb{Z}}$ s.t. $\text{Prob}(X_i = \xi) = p_\xi$ ($\xi = 1, \dots, N$)

Defⁿ. Let T be an invertible map on a probability space. We say that T simulates a Bernoulli process $\mathcal{B}(p_1, \dots, p_N)$, if \exists measurable function

$$f: \Omega \rightarrow \{1, \dots, N\}$$

s.t. $\{f \circ T^k\}_{k \in \mathbb{Z}}$ is a Bernoulli process $\mathcal{B}(p_1, \dots, p_N)$

In other words, not knowing ω , we cannot distinguish using the methods of statistics alone the time series $\{f(T^k \omega) : k \in \mathbb{Z}\}$ from a purely random signal.

Recall:

Defⁿ. The metric entropy of a probability preserving map T on a probability space $(\Omega, \mathcal{F}, \mu)$ is

$$h_\mu(T) = \sup \left\{ h_\mu(T, \alpha) : \begin{array}{l} \alpha \text{ finite} \\ \text{measurable partition} \end{array} \right\}$$

Sinai's Factor Thm. Let T be an invertible probability preserving map with entropy $h_\mu(T)$. Let $p = (p_1, \dots, p_N)$ be a probability vector.

(1) If $-\sum_i p_i \ln p_i \leq h_\mu(T)$, then T simulates a Bernoulli process $B(p_1, \dots, p_N)$

(2) If $-\sum_i p_i \ln p_i > h_\mu(T)$, then it doesn't.

In Summary:

$$\left(\begin{array}{c} \text{positive metric} \\ \text{entropy} \end{array} \right) \iff \left(\begin{array}{c} \exists \text{ experiment } \alpha \\ \text{with completely} \\ \text{random looking} \\ \text{time series} \end{array} \right)$$

Take $\mathcal{L} = \{ \{ \omega \in \Omega : f(\omega) = k \} \mid k = 1, \dots, N \}$