

IV.1 INFORMATION CONTENT

4.1 UNCERTAINTY

Let  $\bullet \xi = \{A_1, A_2, \dots, A_N\}$  be a measurable partition of  $(X, \mathcal{B}, \mu)$   
 $\bullet T: X \rightarrow X$  measure preserving

Let  $\xi(x) = i : \Leftrightarrow x \in A_i$

Idea:  $\xi(x)$  is all we can observe of the state.

"itinerary"

Q: Given  $\xi(x), \xi(Tx), \dots, \xi(T^{n-1}x)$ , how much uncertainty do we have regarding  $\xi(T^n x)$ ? Or: how much information content do we gain by learning  $\xi(T^n x)$ ?

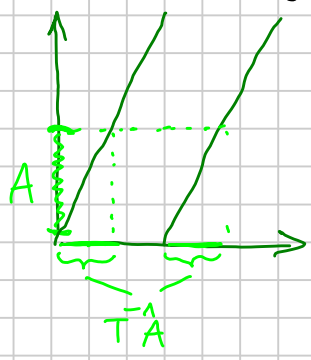
Examples :

1)  $X = S^1, \mu = \text{Lebesgue}, T(x) = x + \frac{1}{100} \pmod{1}$ .  $\xi = \{A_1, A_2\}, A_1 = [0, \frac{1}{2}), A_2 = [\frac{1}{2}, 1)$

Itinerary: 1, 1, 1, 1, 1  $\rightarrow$  If  $x \sim \mu$ , then the next partition element is  $A_1$  with  $\frac{45}{46} \approx 97\%$  probability

2) As 1), only now  $T(x) = 2x \pmod{1}$ .

$\rightarrow$  Here, after any observed itinerary there is a 50% chance of observing both  $A_1$  and  $A_2$  as the next partition element



We would like to be able to describe the difference between examples 1) and 2) quantitatively.

## 4.2 MEASURABLE PARTITION

58

In the entire chapter:  $(X, \mathcal{B}, \mu)$  — prob. space

- A (finite measurable) partition of  $(X, \mathcal{B}, \mu)$ , or, in shorthand, of  $X$ , is a finite collection  $\xi = \{A_1, \dots, A_m\}$  of m'ble sets  $A_i \in \mathcal{B}$ , s.t.

(i) the  $A_i$  are essentially disjoint:  $\mu(A_i \cap A_j) = 0, i \neq j$ ;

(ii)  $\xi$  is a covering:  $\mu\left(\bigcup_{i=1}^m A_i\right) = 1$

- $\xi'$  is a refinement of  $\xi$ , denoted by  $\xi \leq \xi'$ , if  $\xi \subseteq \xi' \pmod{0}$ , i.e. up to differences of measure zero

- The common refinement of partitions  $\xi$  and  $\eta$  is defined as

$$\xi \vee \eta = \{A \cap B \mid A \in \xi, B \in \eta\}$$

- $\xi$  and  $\eta$  are called independent, denoted by  $\xi \perp \eta$ , if

$$\mu(A \cap B) = \mu(A)\mu(B) \quad \forall A \in \xi, B \in \eta$$

Note the trivial one-to-one correspondence between finite partitions and finite  $\sigma$ -algebras.

## 4.3 ENTROPY OF A PARTITION

Suppose  $Y \sim \mu$  is a random variable.

Q: How much information do we gain by learning  $Y \in A$  for  $A \in \mathcal{B}$ ?

↳ Call this  $I(A)$

Natural requirements:

(i)  $I(A) \geq 0$ , and decreasing in  $A$ , i.e.  $I(A) \geq I(B)$  if  $A \subseteq B$ ; and  $I(A) = 0$  if  $\mu(A) = 1$

(ii)  $I$  continuous in  $A$ , i.e.  $I(A_\varepsilon) \rightarrow I(A)$  for  $A_\varepsilon \downarrow A$  or  $A_\varepsilon \uparrow A$

(iii) If  $A, B$  independent, i.e.  $\mu(A \cap B) = \mu(A)\mu(B)$ , then  $I(A \cap B) = I(A) + I(B)$

It turns out, the only function  $\varphi: [0,1] \rightarrow \mathbb{R}_{\geq 0}$  such that  $I(A) = \varphi(\mu(A))$  satisfies (i)-(iii) is  $\varphi(t) = -c \log t$ ,  $c \geq 0$ .

Def: (a) The **information content** of  $A \in \mathcal{B}$  is  $I_{\mu}(A) := -\log \mu(A)$   
(b) The **information function** of a finite measurable partition  $\xi$  is

$$I_{\mu}(\xi)(x) := \sum_{A \in \xi} I_{\mu}(A) \chi_A(x) = - \sum_{A \in \xi} \log \mu(A) \chi_A(x)$$

(c) The **entropy** of a finite measurable partition  $\xi$  is its average information content:

$$H_{\mu}(\xi) := \int I_{\mu}(\xi) d\mu = - \sum_{A \in \xi} \mu(A) \log \mu(A)$$

Convention:  
 $0 \cdot \log 0 = 0$

If  $\mu$  is unambiguous, we will write  $I(A)$ ,  $I(\xi)$ , and  $H(\xi)$ .

Prop: Let  $\xi, \eta$  be finite partitions. Then

- (a)  $H(\xi) \geq 0$ , and  $H(\xi) = 0$  if and only if  $\xi = \{X\}$
- (b) If  $\xi \leq \eta$ , then  $H(\xi) \leq H(\eta)$  with equality only if  $\xi = \eta$
- (c) If  $|\xi| = m$  ( $m$  elements in  $\xi$ ), then  $H(\xi) \leq \log m$ , equality only if  $\mu(A) = 1/m$  for every  $A \in \xi$ .
- (d)  $H(\xi \vee \eta) \leq H(\xi) + H(\eta)$  with equality only if  $\xi \perp \eta$ .

Useful Lemma: Let  $\varphi(t) := -t \log t$ , then for every prob. vector  $(p_1, \dots, p_m)$  and  $x_1, \dots, x_m \in [0,1]$  one has

$$\varphi\left(\sum_i p_i x_i\right) \geq \sum_i p_i \varphi(x_i),$$

with equality only if all the  $x_i$  where  $p_i \neq 0$  are equal.

Pf: (a)-(c): Exercise.

(d) Shorthands:  $\xi = \{A_1, \dots, A_m\}$ ,  $\eta = \{B_1, \dots, B_m\}$   
 $\mu_i = \mu(A_i)$ ,  $\nu_j = \mu(B_j)$ ,  $\kappa_{ij} = \mu(A_i \cap B_j)$ ,  
 hence  $\sum_i \kappa_{ij} = \nu_j$ ,  $\sum_j \kappa_{ij} = \mu_i$

Lemma (with  $x_i = \frac{k_{ij}}{\mu_i}$ ,  $p_i = \mu_i$ ) gives

$$-v_j \log v_j \geq -\sum_i \mu_i \frac{k_{ij}}{\mu_i} \log \frac{k_{ij}}{\mu_i} = \sum_i k_{ij} \log \mu_i - \sum_i k_{ij} \log k_{ij}$$

Sum up over  $j$ :

$$H(\eta) = -\sum_j v_j \log v_j \geq \sum_i \mu_i \log \mu_i - \sum_{ij} k_{ij} \log k_{ij}$$

$$= -H(\xi) + H(\xi \vee \eta)$$

Equality only if for fixed  $j$ , all the  $x_i = \frac{k_{ij}}{\mu_i}$  are the same (without loss:  $\mu_i \neq 0 \forall i$ ), say  $\frac{k_{ij}}{\mu_i} = c_j$

$$\stackrel{\sum_i}{\implies} \sum_i k_{ij} = c_j \sum_i \mu_i = c_j \implies k_{ij} = \mu_i v_j \quad \forall ij, \text{ which is exactly } \xi \perp \eta. \quad \blacksquare$$

## 4.4 CONDITIONAL ENTROPY

$\xi, \eta$ : finite partitions

Q: How much information do we gain by learning  $Y \in A \in \xi$ , if we know  $Y \in B \in \eta$ ?

We already know  $Y \in B \implies$  new probability space  $(\mathcal{B}, \mathcal{B}|_B, \mu_B)$ , where  $\mu_B(E) = \frac{\mu(E \cap B)}{\mu(B)}$ . New information  $Y \in A$  tells us  $Y \in A \cap B$ , hence the information content:

$$I_\mu(A|B) = -\log \frac{\mu(A \cap B)}{\mu(B)} =: -\log \mu(A|B)$$

(content)

If we are not told  $B \in \eta$  yet, just  $\eta$ , the information gain by learning  $Y \in A$  is going to depend on  $Y$  itself:

$$I_\mu(A|\eta)(x) = \sum_{B \in \eta} I_\mu(A|B) \chi_B(x) = -\sum_{B \in \eta} \log \mu(A|B) \chi_B(x) = -\log E[\chi_A|\eta](x)$$

conditional expectation; recall the analogy between  $\eta$  and  $\sigma(\eta)$

If, in addition, we are only told  $\xi, \eta$ , the information content 61  
 is further going to depend on  $\xi(x)$ :

$$I_\mu(\xi|\eta) = \sum_{A \in \xi} I_\mu(A|\eta) \chi_A$$

"(conditional) information function"

This motivates:

Def: The conditional entropy of  $\xi$  given  $\eta$  is the average information gain:

$$H_\mu(\xi|\eta) = \int I_\mu(\xi|\eta) d\mu$$

If  $\xi(x) = A \Leftrightarrow x \in A \in \xi$ , we can write

$$H(\xi|\eta) = \int -\log \mu(\xi(x)|\eta(x)) d\mu(x)$$

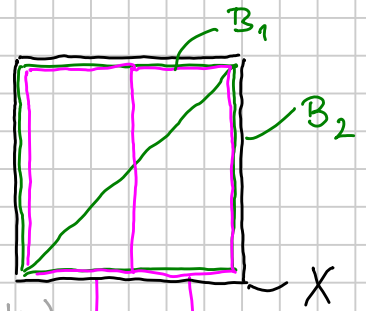
Prop: We have

$$H(\xi \vee \eta) = H(\eta) + H(\xi|\eta)$$

Proof:

$$\begin{aligned} H(\xi|\eta) &= \int -\log \mu(\xi(x)|\eta(x)) d\mu(x) \\ &= \int -\log \mu(\xi(x) \cap \eta(x)) + \log \mu(\eta(x)) d\mu(x) \\ &= \underbrace{\int -\log \mu(\xi(x) \cap \eta(x)) d\mu(x)}_{= H(\xi \vee \eta)} - \underbrace{\int -\log \mu(\eta(x)) d\mu(x)}_{= H(\eta)} \end{aligned}$$

$$\begin{aligned} \xi &= \{A_1, A_2\} \\ \eta &= \{B_1, B_2\} \end{aligned}$$



Observe:

$$I(A_1|B_2) > I(A_1|B_1)$$

$$I(A_2|B_1) > I(A_2|B_2)$$

## IV 2 ENTROPY OF MEASURE PRESERVING TRANSFORMATIONS

62

### 4.5 METRIC ENTROPY

Throughout: let  $(X, \mathcal{B}, \mu, T)$  be a ppt.

Q: What is the information content in " $x \in A_0, Tx \in A_1, \dots, T^{m-1}x \in A_{m-1}$ "?

If  $A_0, \dots, A_{m-1} \in \mathcal{E}$  then this is the same information as

$$x \in A_0 \cap T^{-1}A_1 \cap \dots \cap T^{-(m-1)}A_{m-1}$$

→ The dynamics induces a partition

Let

•  $T^{-k}\mathcal{E} := \{T^{-k}A \mid A \in \mathcal{E}\}$  — this is a measurable partition

•  $\mathcal{E}^{(m)} := \bigvee_{k=0}^{m-1} T^{-k}\mathcal{E} = \mathcal{E} \vee T^{-1}\mathcal{E} \vee \dots \vee T^{-(m-1)}\mathcal{E}$

Since  $T$  is a ppt, we have  $h(T^{-k}\mathcal{E}) = h(\mathcal{E}) \quad \forall k \geq 0$

From Prop 4.3:  $h(\mathcal{E} \vee \eta) \leq h(\mathcal{E}) + h(\eta)$

$$\begin{aligned} \Rightarrow h(\mathcal{E}^{(m+m)}) &= h(T^{-m}\mathcal{E}^{(m)} \vee \mathcal{E}^{(m)}) \\ &\leq h(\mathcal{E}^{(m)}) + h(\mathcal{E}^{(m)}) \end{aligned}$$

Thus, the sequence  $\{h(\mathcal{E}^{(m)})\}_{m \in \mathbb{N}}$  is *sub-additive*.

Remark:  $f(0) \geq 0$   
concave  $\Rightarrow$  sub-additive  
 $\nleftarrow$   
" $\rightarrow$ " draw or compute  
" $\nleftarrow$ " take  $|\sin(x)|$

Lemma: For a sub-additive sequence  $\{a_m\}_{m \in \mathbb{N}}$ , i.e. where  $a_{n+m} \leq a_n + a_m$   
 $\forall m, n \in \mathbb{N}$ , one has

$$\lim_{m \rightarrow \infty} \frac{1}{m} a_m = \inf_m \frac{a_m}{m}$$

Pf: Fix  $n$ . For arbitrary  $m \in \mathbb{N}$ , write  $m = km + r$  with  $r \in \{0, \dots, m-1\}$

Then  $a_m \leq k a_n + a_r$  by sub-additivity, thus

$$\frac{a_m}{m} \leq \frac{k a_n + a_r}{km + r} \leq \frac{a_n}{m} + \frac{a_r}{m}$$

$$\Rightarrow \limsup_{m \rightarrow \infty} \frac{a_m}{m} \leq \frac{a_n}{m} \quad \forall n \Rightarrow \limsup_{m \rightarrow \infty} \frac{a_m}{m} \leq \inf_n \frac{a_n}{m}$$

$$\text{Also, clearly, } \liminf_{m \rightarrow \infty} \frac{a_m}{m} \geq \inf_n \frac{a_n}{n}$$

$$\implies \lim_{m \rightarrow \infty} \frac{a_m}{m} = \inf_n \frac{a_n}{n}$$

By this, it makes sense to speak of the average information content generated by  $T$  in one step:

Def: The metric (or measure-theoretic) entropy of  $T$  wrt the partition  $\xi$  is defined as

$$h(T, \xi) := \lim_{m \rightarrow \infty} \frac{1}{m} H(\xi^{(m)})$$

(Supposed  $H(\xi) < \infty$ .)

Recall Q: How much information does  $\xi(T^m x)$  contain, given  $\xi(x), \dots, \xi(T^{m-1}x)$ ?

Prop [BS, Prop 9.3.1]:

$$h(T, \xi) = \lim_{m \rightarrow \infty} H(\xi | T^{-1} \xi^{(m)}) = \lim_{m \rightarrow \infty} H(\xi | \bigvee_{i=1}^m T^{-i} \xi)$$

Proof: It holds  $H(\xi | \zeta) \leq H(\xi | \eta)$  if  $\eta \leq \zeta$ , where  $\xi, \eta, \zeta$  are finite partitions (exercise). Hence,  $H(\xi | T^{-1} \xi^{(m)})$  is a non-increasing sequence in  $m$ .

$$H(\xi^{(m)}) = H(\xi \vee T^{-1} \xi^{(m-1)})$$

$$\stackrel{\text{Prop 4.4}}{=} H(T^{-1} \xi^{(m-1)}) + H(\xi | T^{-1} \xi^{(m-1)})$$

$$\stackrel{T \text{ ppt} \rightarrow H(T^{-1} \eta) = H(\eta)}{=} H(\xi^{(m-1)}) + H(\xi | T^{-1} \xi^{(m-1)})$$

$$= \dots = H(\xi) + \sum_{k=1}^{m-1} H(\xi | T^{-k} \xi)$$

Divide by  $m$ , and let  $m \rightarrow \infty$  to obtain the claim

Def: The metric (or measure-theoretic) entropy of  $T$  is defined as [64]

$$h(T) = \sup_{\xi} h(T, \xi),$$

where the sup is over all finite measurable partitions.

#### 4.6 KOLMOGOROV-SINAI GENERATOR THEOREM

$d(\xi, \eta) := \min_{\sigma} \sum_{i=1}^m \mu(A_i \Delta B_{\sigma(i)})$  is a metric for finite measurable partitions  
 permutation  $\uparrow$   $\xi = \{A_1, \dots, A_m\}, \eta = \{B_1, \dots, B_m, \underbrace{\phi, \dots, \phi}_{\text{add empty sets until } |\xi|=|\eta|}\}$

Def: (a) The sequence  $\{\xi_n\}_m$  of finite partitions is **refining** if  $\xi_m \leq \xi_{m+1}$   $\forall m$

(b) The sequence  $\{\xi_n\}_m$  of finite partitions is **generating** if for every finite partition  $\eta$  and  $\delta > 0$  there is a  $\tilde{\eta} \leq \bigvee_{i=1}^m \xi_i$  for a sufficiently large  $m$  such that  $d(\eta, \tilde{\eta}) < \delta$ .

(c) A **generator** for a ppt  $T$  is a finite partition  $\xi$  such that  $\{\xi^{(m)}\}_m$  is generating.

Note: Part (b) is a quantitative version of  $\sigma\left(\bigcup_{m=1}^{\infty} \xi_m\right) = \mathcal{B} \text{ mod } \mu$

Example:  $\xi = \left\{ \left[0, \frac{1}{2}\right), \left[\frac{1}{2}, 1\right) \right\}, T(x) = 2x \text{ mod } 1.$

$\xi$  is a generator for  $T$ , since  $\bigcup_{n=0}^{\infty} \bigvee_{i=0}^m T^{-i} \xi$  contains all dyadic intervals, hence  $\sigma\left(\bigcup_{m=1}^{\infty} \xi^{(m)}\right) = \mathcal{B}$ , the Borel  $\sigma$ -algebra (cf. §1.17)

Prop: If  $\{\xi_n\}_m$  is a refining and generating sequence of partitions, then  $h(T) = \lim_{m \rightarrow \infty} h(T, \xi_m)$ . [BS, Prop. 9.3.4]

Theorem (Kolmogorov-Sinai generator theorem):

Let  $\xi$  be a generator for  $T$ . Then  $h(T) = h(T, \xi)$ . [BS, Thm 9.3.5]

#### 4.7 THE SHANNON-MCMILLAN-BREIMAN THEOREM

Theorem: Let  $T$  be an ergodic ppt and  $\xi$  a finite partition. Then

$$\frac{1}{n} \sum_{i=0}^{n-1} \log \mu(\xi^{(i)}) \xrightarrow[n \rightarrow \infty]{} h(T, \xi) \quad \mu\text{-a.e. and in } L^1$$

[Sa, Thm. 4.3]



Consequence: Recalling Def 4.3 (information function), the theorem 1.65  
states that at a typical point  $x \in X$ , the size of the partition element  $\xi^{(m)}(x)$ ,  
 $\mu(\xi^{(m)}(x))$  decays as  $\exp(-m h(T, \xi))$

## 4.8 EXAMPLES

1) Circle rotation:  $T(x) = x + d \pmod{1}$  on  $S^1$

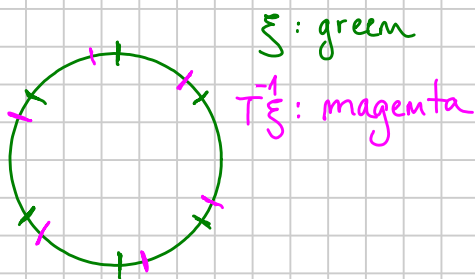
(i)  $d \in \mathbb{Q}$ :  $\exists k \in \mathbb{N}$  s.t.  $T^{-k} = \text{id}$ , hence  $T^{-k}\xi = \xi$  for every partition  $\xi$ ,  
and thus  $\xi \leq \bigvee_{i=1}^m T^{-i}\xi \quad \forall m \geq k$ .

But from Prop 4.4:  $H(\xi \vee \eta) = H(\eta) + H(\xi|\eta)$ , and if  
 $\xi \leq \eta$ , then  $H(\xi|\eta) = 0$  follows.

$$\Rightarrow H(\xi | \bigvee_{i=1}^m T^{-i}\xi) = 0 \quad \forall m \geq k \xrightarrow{\text{Prop 4.5}} h(T, \xi) = 0$$

$$\Rightarrow h(T) = \sup_{\xi} h(T, \xi) = 0$$

(ii)  $d \notin \mathbb{Q}$ :  $\xi_N$ : partition of  $S^1$  into  $N$  equal intervals



Since  $d \notin \mathbb{Q}$ , endpoints of  $T^{-i}\xi$   
are mutually distinct

$\Rightarrow \xi^{(m)}$  consists of  $m \cdot N$   
intervals, and the endpoints  
of  $\xi^{(0)}, \xi^{(1)}, \dots$  are dense in  $S^1$

$\Rightarrow \xi_N$  is generator (every interval can be arbitrary  
well approximated)

$$\text{Prop 4.3 (c): } H(\xi_N^{(m)}) \leq \log(mN) = \log m + \log N$$

$$\Rightarrow h(T) = h(T, \xi_N) = 0$$

↑ Kolmogorov-Sinai thm

2) Expanding maps:  $T(x) = kx \pmod{1}$  on  $S^1$

$$\text{Let } \xi = \left\{ \left[0, \frac{1}{k}\right), \left[\frac{1}{k}, \frac{2}{k}\right), \dots, \left[\frac{k-1}{k}, 1\right) \right\}$$

Easy to see that

$$\xi^{(m)} = \left\{ \left[ \frac{i}{2^m}, \frac{i+1}{2^m} \right) \mid i = 0, \dots, 2^m - 1 \right\}$$

and thus  $\xi$  is a generator for  $T$ .

$$H(\xi^{(m)}) = - \sum_{i=1}^{2^m} \frac{1}{2^m} \log \left( \frac{1}{2^m} \right) = m \log 2 \xrightarrow{\text{Kolmogorov-Sinai}} h(T) = \log 2$$

3)  $\left(\frac{1}{2}, \frac{1}{2}\right)$ -Bernoulli shift:

Recall cylinders:  $[a_0, \dots, a_k] = \left\{ \{x_i\}_{i \in \mathbb{N}} \mid x_j = a_j, j = 0, \dots, k \right\}$

$\xi = \{[0], [1]\}$  is generating, since  $T^{-i}[a] = \left[ \underbrace{*}_{i}, \dots, *, a \right]$ , hence

$\bigvee_{i=0}^k T^{-i} \xi$  contains exactly the cylinders  $[a_0, \dots, a_k]$ ,  $a_i \in \{0, 1\}, i = 0, \dots, k$

Thus

$$\begin{aligned} H(\xi^{(m)}) &= - \sum_{a_0, \dots, a_m \in \{0, 1\}} \mu[a_0, \dots, a_m] \cdot \log \mu[a_0, \dots, a_m] \\ &= \sum_{i=1}^{2^{m+1}} \frac{1}{2^{m+1}} \log(2^{m+1}) = (m+1) \log 2 \end{aligned}$$

$$\Rightarrow h(T) = \log 2$$

4)  $(P, P)$ -Markov shift:

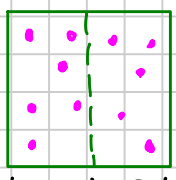
Similar computation as in 3) (cf [Wa, Thm 4.27], [BS, Sec 9.4], [Sa, Prop. 4.7])

yields  $h(T) = - \sum_{i,j} p_i p_{ij} \log p_{ij}$

4.9 HISTORICAL REMARKS

- Boltzmann (1872) / Gibbs (1878): statistical mechanics & thermodynamics

$N$  particles



microstate: each particle left/right

macrostate: # particles left/right

Q: Given macrostate, how much uncertainty about microstate?

Macro:  $(m, N-m) \Rightarrow$  # possible micro =  $\binom{N}{m} \rightsquigarrow$  max if  $m = N/2 \rightsquigarrow$  pressure equilibrium

• von Neumann (~1940): isomorphy of mpts

Q: Are the  $(\frac{1}{2}, \frac{1}{2})$ -Bernoulli and  $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ -Bernoulli shifts isomorph?

• Shannon (1948): entropy & information theory

• Kolmogorov, Sinai (1958): entropy of mpts

Thm: Entropy is invariant under measure-theoretic isomorphism

A:  $\log 2 \neq \log 3$ , hence they are not isomorph

• Ornstein (1970): Bernoulli shifts with same entropy are isomorph

4.10 TOPOLOGICAL ENTROPY

•  $X$ : compact metric space

•  $T: X \rightarrow X$  continuous

•  $\xi$ : open cover of  $X$

$N(\xi) := \min \{ |\eta| \mid \eta \subseteq \xi, X \subseteq \bigcup_{U \in \eta} U \}$  smallest cardinality of finite subcover

It holds  $N(\xi \vee \eta) \leq N(\xi) \cdot N(\eta)$  (\*)

Q: At which rate is  $T$  refining  $\xi$ ?  $\rightarrow$  Look at  $N\left(\bigvee_{k=0}^{m-1} T^{-k}\xi\right)$ ,  $m \geq 0$

By (\*), growth is subexponential in  $m$ .

Topological entropy

• of  $T$  relative to  $\xi$ :  $h_{top}(T, \xi) = \lim_{m \rightarrow \infty} \frac{1}{m} \log N\left(\bigvee_{k=0}^{m-1} T^{-k}\xi\right)$

• of  $T$ :  $h_{top}(T) = \sup_{\xi} h_{top}(T, \xi)$   
 $\xi$  finite open covers of  $X$

Theorem (Variational principle)

$T: X \rightarrow X$  continuous on the compact metric space  $X$ . Then

$h_{top}(T) = \sup \{ h_{\mu}(T) \mid \mu \text{ is invariant Borel measure} \}$