Entropy in Ergodic Theory

Tomasz Downarowicz

Institute of Mathematics and Computer Science Wroclaw University of Technology Poland

Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics Janu

January 28–February 1, 2013

1/24

Based on

ENTROPY IN DYNAMICAL SYSTEMS

New Mathematical Monographs: 18 Cambridge University Press 2011

PART I: Entropy in Ergodic Theory

Tomasz Downarowicz (Poland)

2



Tomasz Downarowicz (Poland)

January 28–February 1, 2013 3/24

<ロ> <同> <同> <同> <同> < 同>



Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics

January 28–February 1, 2013 3

<ロ> <同> <同> <同> <同> < 同>

3/24



Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics

January 28–February 1, 2013 3

3/24

э



Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics

January 28–February 1, 2013 3

<ロ> <同> <同> <同> <同> < 同>

3/24

How much information was that?

January 28–February 1, 2013 4 / 24

э

How much information was that?

one out of two choices = ONE BIT

э

イロト イポト イヨト イヨト



Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics Janu

January 28–February 1, 2013 5

<ロ> <同> <同> <同> <同> < 同>

5/24



Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics Janu

January 28–February 1, 2013 5 /

5/24



Tomasz Downarowicz (Poland)

January 28–February 1, 2013 5

5/24



one of four choices = TWO BITS

Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics

2

・ロト ・ 四ト ・ ヨト ・ ヨト



Tomasz Downarowicz (Poland)

January 28-February 1, 2013 5/24

イロト イロト イヨト イヨト



one of three choices = ONE AND HALF BITS

Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics

January 28–February 1, 2013 5 / 24

э

・ロト ・ 四ト ・ ヨト ・ ヨト



NO - this SCHOOL is about NONLINEAR SCIENCE !!!

Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics

A B b 4 B b January 28-February 1, 2013 5/24

< 47 ▶

э

0 BITS = 1 choice

1 BIT = 2 choices 2 BITS = 4 choices 3 BITS = 8 choices etc.

BITS = log₂(# choices)

3 choices = $\log_2(3)$ BITS \approx 1.585 BITS





DID I WIN? (YES/NO)

Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics

January 28–February 1, 2013 7 / 24

イロト イヨト イヨト イヨト

NO - 999999 of a million chances I KNEW IT ANYWAY...

(there was almost only one choice - nearly no information gained)

YES - 1 of a million chances HURRA!!! (large information gained)

YES - 1 of a million chances HURRA!!! (large information gained) # BITS = $log_2(1000000) \approx 19,9$ BITS

NO - 999999 of a million chances I KNEW IT ANYWAY...

(there was almost only one choice - nearly no information gained)

BITS = ???

YES - 1 of a million chances HURRA!!! (large information gained) # BITS = $\log_2(1000000) \approx 19,9$ BITS

 $\log_2(1000000) = -\log_2(\frac{1}{1000000}) = -\log_2(\text{probability of winning})$

NO - 999999 of a million chances I KNEW IT ANYWAY...

(there was almost only one choice - nearly no information gained)

BITS = $-\log_2(\text{probability of loosing}) = -\log_2(\frac{999999}{1000000}) \approx 0,0000014$

DEFINITION 1

If Ω is a finite probability space with atoms x_1, x_2, \ldots of probabilities $P(x_i)$, $(i = 1, 2, \ldots)$, then the associated *information function* on Ω is defined as

 $I(x_i) = -\log_2(P(x_i)).$

DEFINITION 1

If Ω is a finite probability space with atoms x_1, x_2, \ldots of probabilities $P(x_i)$, $(i = 1, 2, \ldots)$, then the associated *information function* on Ω is defined as

$$I(x_i) = -\log_2(P(x_i)).$$

If Ω is finite and has *n* elements of equal probabilities $\frac{1}{n}$ then the information function function is constant equal everywhere to $\log_2(n)$.

DEFINITION 2

If (Ω, Σ, μ) is a (perhaps non-atomic) probability space and $\mathcal{P} = \{P_1, P_2, \dots\}$ is a countable (or finite) measurable partition of Ω then the associated *information function* on Ω is defined as

$$I_{\mathcal{P}}(\mathbf{x}) = -\log_2(\mu(\mathbf{P}_{\mathbf{x}})),$$

where P_x is the unique element of \mathcal{P} such that $P_x \ni x$.



Tomasz Downarowicz (Poland)

イロト イロト イヨト イヨト

11/24





・ロト ・ 四ト ・ ヨト ・ ヨト



・ロト ・ 四ト ・ ヨト ・ ヨト



Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics January 28–February 1, 2013 11 / 24



<ロ> <四> <四> <四> <四> <四</p>

DEFINITION 3

If (Ω, Σ, μ) is a probability space and $\mathcal{P} = \{P_1, P_2, ...\}$ is a countable measurable partition of Ω then the *Shannon entropy* of \mathcal{P} is defined as the expected value of the information function:

$$H(\mathcal{P}) = \int I_{\mathcal{P}} d\mu = -\sum_{i} \mu(P_{i}) \log_{2} \mu(P_{i})$$

(The average over the space information delivered by the partition.)

イロト 不得 トイヨト イヨト



Consider the two bitmaps





э

Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics January 28–February 1, 2013 13 / 24



Consider the two bitmaps



They have the same sizes (even the same proportion of black and white). Thus they carry the same Shannon information (= # pixels). However...

Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics

13/24


Consider the two bitmaps



Any zipping program compresses the left hand side bitmap about 5 times more than the right hand side bitmap. Why?

Tomasz Downarowicz (Poland)



Consider the two bitmaps



Imagine that you explain how to draw each bitmap over the phone... How much INFORMATION is needed for each of them?

What makes the difference between these bitmaps, if both carry the same Shannon information?

イロト イポト イヨト イヨト

What makes the difference between these bitmaps, if both carry the same Shannon information?

The answer is delivered by the dynamic entropy and the Shannon–McMillan–Breiman Theorem.

イロト 不得 トイヨト イヨト

Now we will assume that on our probability space (Ω, Σ, μ) we have a measurable transformation $T : \Omega \to \Omega$ which preserves the measure μ , that is $\mu(T^{-1}(A)) = \mu(A)$ for every $A \in \Sigma$.

(日)

Now we will assume that on our probability space (Ω, Σ, μ) we have a measurable transformation $T : \Omega \to \Omega$ which preserves the measure μ , that is $\mu(T^{-1}(A)) = \mu(A)$ for every $A \in \Sigma$.

EXAMPLE

Let $\Omega = \{0, 1\}^{\mathbb{N}}$, $T = \text{shift} (T(x_1, x_2, ...) = (x_2, x_3, ...))$ and μ is some shift-invariant measure. Every such measure is determined by its values on cylinders $C = [c_1, c_2, ..., c_n]$.

< 白 > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >



Tomasz Downarowicz (Poland)

1, 2013 16 / 24



イロト イロト イヨト イヨト

э



Tomasz Downarowicz (Poland)



Tomasz Downarowicz (Poland)



Tomasz Downarowicz (Poland)

・ロト ・ 四ト ・ ヨト ・ ヨト

э



<ロト < 回 > < 回 > < 回 > < 回 > … 回









・ロト ・ 四ト ・ ヨト ・ ヨト



Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics January 28–February 1, 2013 16 / 24

・ロト ・ 四ト ・ ヨト ・ ヨト

æ



DEFINITION 4

Let (Ω, Σ, μ) be a probability space and let $T : \Omega \to \Omega$ be a measurable and measure-preserving transformation. Let $\mathcal{P} = \{P_1, P_2, ...\}$ be a countable measurable partition of Ω . Then the *information function in n steps* on Ω is defined as

$$I_{\mathcal{P}^n}(\mathbf{x}) = -\log_2(\mu(\mathbf{P}_{\mathbf{x}}^n)),$$

where

$$P_{x}^{n} = P_{x} \cap T^{-1}(P_{Tx}) \cap T^{-2}(P_{T^{2}x}) \cap \cdots \cap T^{-n+1}(P_{T^{n-1}x})$$

(it is the unique element of the partition $\mathcal{P}^n := \bigvee_{i=0}^{n-1} T^{-i}(\mathcal{P})$ containing *x*, and is called the *n*-cylinder of *x*).

イロト 不得 トイヨト イヨト



$x \in [0,1]$

Tomasz Downarowicz (Poland) Recent Trends in Nonlinear Dynamics January 28–February 1, 2013

▶ ≣ ∽۹<0 2013 18/24

ヘロト ヘロト ヘヨト ヘヨト



(Poland) Recent Trends in Nonlinear Dynamics January 28–February 1, 2013 18 / 24

ヘロト ヘロト ヘヨト ヘヨト

-2

Tomasz Downarowicz (Poland)



$$x \in [0, 1]$$

 $x = 0.765900862...$

To fully identify *x* we need infinite amount of information.

Recent Trends in Nonlinear Dynamics January 28–February 1, 2013 18 / 24

イロン イ理 とく ヨン イヨン

э



To fully identify x we need infinite amount of information.

With each digit we acquire $log_2(10)$ BITS of information.



To fully identify *x* we need infinite amount of information.

With each digit we acquire $log_2(10)$ BITS of information.

This corresponds to the flow of information in the dynamical system:



To fully identify *x* we need infinite amount of information.

With each digit we acquire $log_2(10)$ BITS of information.

This corresponds to the flow of information in the dynamical system:

 $T: [0, 1] \rightarrow [0, 1],$ $T(x) = 10x \mod 1,$ μ is the Lebesgue measure



To fully identify *x* we need infinite amount of information.

With each digit we acquire $log_2(10)$ BITS of information.

This corresponds to the flow of information in the dynamical system:

 $T: [0, 1] \rightarrow [0, 1],$ $T(x) = 10x \mod 1,$ μ is the Lebesgue measure T(0.765900862...) = 0.65900862..., $T^{2}(0.765900862...) = T(0.65900862...) = 0.5900862...$

A D A A B A A B A A B A B B



To fully identify *x* we need infinite amount of information.

With each digit we acquire $log_2(10)$ BITS of information.

This corresponds to the flow of information in the dynamical system:

 $\begin{array}{l} T:[0,1] \rightarrow [0,1], \\ T(x) = 10x \ \mathrm{mod} \ 1, \\ \mu \ \mathrm{is \ the \ Lebesgue \ measure} \\ T(0.765900862...) = 0.65900862..., \\ T^2(0.765900862...) = T(0.65900862...) = 0.5900862... \end{array}$

 $\mathcal{P} = \{[0, 0.1), [0.1, 0.2), \dots, [0.9, 1]\}$

• $P^n(x) =$

{points that give the same answers as *x* through *n* times}

A D F A B F A B F A B F

- *Pⁿ(x)* = {points that give the same answers as *x* through *n* times}
- $I_{\mathcal{P}^n}(\mathbf{x}) = -\log \mu(\mathcal{P}^n(\mathbf{x}))$

- *Pⁿ(x)* = {points that give the same answers as *x* through *n* times}
- $I_{\mathcal{P}^n}(\mathbf{x}) = -\log \mu(\mathbf{P}^n(\mathbf{x}))$
- $H(\mathcal{P}^n) := \int I_{\mathcal{P}^n} d\mu$ (average *over space* information in *n* steps)

(日)

- *Pⁿ(x)* = {points that give the same answers as *x* through *n* times}
- $I_{\mathcal{P}^n}(\mathbf{x}) = -\log \mu(\mathbf{P}^n(\mathbf{x}))$
- $H(\mathcal{P}^n) := \int I_{\mathcal{P}^n} d\mu$ (average over space information in *n* steps)

DEFINITION 5

The dynamic entropy of the partition \mathcal{P} is defined as

$$h(T,\mathcal{P}):=\lim_{n}\frac{1}{n}H(\mathcal{P}^{n}).$$

- *Pⁿ(x)* = {points that give the same answers as *x* through *n* times}
- $I_{\mathcal{P}^n}(\mathbf{x}) = -\log \mu(\mathbf{P}^n(\mathbf{x}))$
- $H(\mathcal{P}^n) := \int I_{\mathcal{P}^n} d\mu$ (average over space information in *n* steps)

DEFINITION 5

The dynamic entropy of the partition \mathcal{P} is defined as

$$h(T,\mathcal{P}) := \lim_{n} \frac{1}{n} H(\mathcal{P}^{n}).$$

The dynamic entropy is interpreted as the average over space and *time* gain of information per step.

(日)

Shannon–McMillan–Breiman Theorem

Tomasz Downarowicz (Poland) Recent Trends in Nonlinear Dynamics January 28–February 1, 2013 20 / 24

・ロト ・ 四ト ・ ヨト ・ ヨト

э

Shannon–McMillan–Breiman Theorem

THEOREM 1

If μ ergodic then

$$\frac{1}{n} l_{\mathcal{P}^n}(\mathbf{x}) \xrightarrow[n \to \infty]{\mu-a.e.} h(T, \mathcal{P})$$

Tomasz Downarowicz (Poland)

Recent Trends in Nonlinear Dynamics January 28–February 1, 2013

20/24

Shannon–McMillan–Breiman Theorem

THEOREM 1

If μ ergodic then

$$\frac{1}{n} I_{\mathcal{P}^n}(\mathbf{x}) \stackrel{\mu-a.e.}{\underset{n\to\infty}{\longrightarrow}} h(T,\mathcal{P})$$

That is, the *average gain of information per step* does not depend on the initial point.
Let $\Omega = \{0, 1\}^{\mathbb{N}}$, T = shift and μ is some ergodic shift-invariant measure. Then for a μ -"typical" point $x = (x_1, x_2, ...)$ the measure of a long initial cylinder $x[1, n] := [x_1, x_2, ..., x_n]$ is approximately $2^{-nh(T, \mathcal{P})}$, where \mathcal{P} is the two-element partition $\{[0], [1]\}$.

Let $\Omega = \{0, 1\}^{\mathbb{N}}$, T = shift and μ is some ergodic shift-invariant measure. Then for a μ -"typical" point $x = (x_1, x_2, ...)$ the measure of a long initial cylinder $x[1, n] := [x_1, x_2, ..., x_n]$ is approximately $2^{-nh(T, \mathcal{P})}$, where \mathcal{P} is the two-element partition $\{[0], [1]\}$.

The meaning of "approximately" is very rough, it means only that $-\frac{1}{n}\log_2 \mu(x[1,n]) \approx h(T, \mathcal{P}).$



• Let us go back to our example with the two bitmaps:





The first bitmap is "highly organized" (in fact periodic), hence has small entropy, the second one is "highly random", hence has large entropy, thus $h_1 \ll h_2$.

The first bitmap is "highly organized" (in fact periodic), hence has small entropy, the second one is "highly random", hence has large entropy, thus $h_1 \ll h_2$.

The entropies represent the *average information contents per symbol*. By the Shannon–McMillan–Breiman Theorem, the same average information contents per symbol occurs already in these orbits (bitmaps).

The first bitmap is "highly organized" (in fact periodic), hence has small entropy, the second one is "highly random", hence has large entropy, thus $h_1 \ll h_2$.

The entropies represent the *average information contents per symbol*. By the Shannon–McMillan–Breiman Theorem, the same average information contents per symbol occurs already in these orbits (bitmaps).

So the *effective information* carried by the bitmaps is proportional to h_1 and h_2 , respectively (times the # of pixels). This explains the huge difference.

イロン イロン イヨン イヨン 三日

Everything that was said in this presentation will be given rigorous explanation during the rest of the course...

Everything that was said in this presentation will be given rigorous explanation during the rest of the course...

using more traditional media, such as blackboard (or whiteboard) and chalk (or markers).