# **05. Entropy in Classical Information Theory 1. Motivation**

- <u>Recall: Gibbs' approach to statistical mechanics</u>
  - $\{\Gamma, \rho\}$  = an ensemble of classical states
  - $\Gamma$  = a phase space of multi-particle microstates x
  - $\rho$  = a Gibbs probability distribution defined on  $\Gamma$
  - $S_{\text{Gibbs}}(\rho)$  = ensemble average of  $-\ln\rho$
- Shannon's approach to classical information
  - Generalize the notion of a classical phase space Γ of microstates x to a random variable X with possible values x.
  - View  $\rho$  as a probability distribution that assigns probabilities to the possible values x of X.
  - View  $-\ln\rho$  as a measure of "information"  $\leftarrow$
- <u>*Ex*</u>: Let  $X = \{x_1, ..., x_\ell\}$  = set of  $\ell$  messages.



The amount of information gained from the reception of a message depends on how *likely* it is. 1. Motivation

- 2.  $S_{\text{Shan}}$  as Uncertainty
- 3.  $S_{\text{Shan}}$  as Message Compression



*Claude Shannon* (1916-2001)

- <u>Intuition</u>: The greater the probability  $\rho(x)$ , the more certain that the value of X is x, and the less information associated with this result.

<u>Intuition</u>: The less likely a message is, the more info gained upon its reception! **Def. 1** (*Shannon entropy*). Let *X* be a random variable with possible values  $\{x_1, ..., x_\ell\}$  and probability distribution  $\{p_1, ..., p_\ell\}$ . The **Shannon entropy**  $S_{\text{Shan}}(X)$  of X is given by

 $S_{\text{Shan}}(X) = -\sum_{i=1}^{\ell} p_i \log_2 p_i$ 

• Compare with *S*<sub>Boltz</sub>:

 $S_{\text{Boltz}}(\Gamma_M) = -Nk \sum_{i=1}^{\ell} p_i \ln p_i + \text{const.} \qquad \longleftarrow \begin{array}{c} p_i \text{ are probabilities defined on} \\ single-particle microstates \end{array}$ 

• Continuous version of *S*<sub>Shan</sub>:

$$S_{\text{Shan}}(X) = -\int_{X} \rho(x) \log_2 \rho(x) dx \quad \checkmark X \text{ takes a continuum of values}$$

• Compare with *S*<sub>Gibbs</sub>:

$$S_{\text{Gibbs}}(\rho) = -k \int_{\Gamma} \rho(x) \ln \rho(x) dx$$
$$\rho(x) \text{ are probabilities defined}$$
on multi-particle microstates

- Why the  $\log_2$  in  $S_{\text{Shan}}$ ?
- Short answer: Classical info is measured in units of "bits".
- $\log_2 x = y$  means  $x = 2^y$

## Why the $\log_2 in S_{\text{Shan}}$ ?

- <u>Long answer</u>:

**Claim** (Shannon 1949).  $S_{\text{Shan}}(X) = -\sum_{i} p_i \log_2 p_i$  is the unique function H(X): {*probability distributions on* X}  $\rightarrow \mathbb{R}$ , that satisfies:

- *Continuity*.  $H(p_1, ..., p_\ell)$  is continuous.
- Additivity.  $H(p_1q_1, ..., p_\ell q_\ell) = H(P) + H(Q)$ , for probability distributions P, Q.
- *Monoticity*. Info increases with  $\ell$  for uniform distributions: If  $m > \ell$ , then H(Q) > H(P), for any  $P = \{1/\ell, ..., 1/\ell\}$  and  $Q = \{1/m, ..., 1/m\}$ .
- *Branching*.  $H(p_1, ..., p_\ell)$  is independent of how the process is divided into parts.
- *Bit normalization*. The average info gained for two equally likely messages is one bit:  $H(\frac{1}{2}, \frac{1}{2}) = 1$ .

Bit renormalization requires log<sub>2</sub>

- <u>Suppose</u>:  $X = \{x_1, x_2\}$ , and  $P = \{\frac{1}{2}, \frac{1}{2}\}$ .
- <u>Then</u>:

- <u>And</u>:  $\log 2 = 1$  if and only if  $\log is$  to base 2.  $\log_2 x = y \Rightarrow x = 2^y$ 

#### Why call this "entropy"?

"Nobody really knows what entropy really is, so in a debate you will always have the advantage."







### 2. S<sub>Shan</sub> as a Measure of Uncertainty

Let X be a random variable with possible values {x<sub>1</sub>, ..., x<sub>ℓ</sub>} and probability distribution {p<sub>1</sub>, ..., p<sub>ℓ</sub>}.

**Def. 3.** The **expected value** E(X) of *X* is given by  $E(X) = \sum_{i=1}^{\ell} p_i x_i$ 

**Def. 4.** The **information gained** if *X* is measured to have the value  $x_i$  is given by  $-\log_2 p_i$ .

• Then the expected value of  $-\log_2 p_i$  is  $S_{\text{Shan}}(X)$ :

$$E(-\log_2 p_i) = -\sum_{i=1}^{\ell} p_i \log_2 p_i = S_{\text{Shan}}(X)$$

- $S_{\text{Shan}}(X)$  is the expected information gained upon measuring X.
- The greater  $S_{\text{Shan}}(X)$ , the greater the info gained upon measuring *X*, and the greater the uncertainty of its measured value.

#### **Uncertainty Interpretation Comparison**

Shannon	Boltzmann	Gibbs
X = random variable	$\Gamma_{\mu} = \text{single-particle phase space}$	$\Gamma$ = multi-particle phase space
$\{x_1,, x_\ell\} = \text{values of } X$	$\{x_1,, x_\ell\} = single-particle$ microstates.	$x \in \Gamma$ : multi-particle microstates.
$\{p_1,, p_\ell\}$ = probabilty distribution over values.	$\{p_1,, p_\ell\} = \text{probabilty}$ distribution on $\Gamma_{\mu}$ .	$ \rho = \text{probabilty distribution} $ on Γ.
$p_i$ = probability that X has value $x_i$ upon measurement.	$p_i$ = probability that microstate $x_i$ of particle is in <i>i</i> th cell of $\Gamma_{\mu}$ .	$\rho(x,t) =$ prob that microstate of system at time <i>t</i> is <i>x</i> .
$-\log_2 p_i = \text{info gained}$ upon measurement of $X$ with outcome $x_i$ .	$-\ln p_i = \text{info gained upon}$ finding a particle to be in microstate $x_i$ in <i>i</i> th cell of $\Gamma_{\mu}$ .	$-\ln\rho(x,t) = \text{info gained upon}$ finding multi-particle system to be in microstate <i>x</i> at time <i>t</i> .
$S_{\rm Shan}(X) = -\sum_i p_i \log_2 p_i$	$S_{\text{Boltz}}(\Gamma_{M_D}) = -Nk\sum_i p_i \ln p_i$	$S_{\rm Gibbs}(\rho) = -\int_{\Gamma} \rho \ln \rho dx$
$S_{\text{Shan}}(X) = \text{expected info gain}$ upon measurement of <i>X</i> .	$S_{\text{Boltz}}/N = \text{expected info gain}$ upon finding a single particle of an <i>N</i> -particle system in microstate $x_i$ in <i>i</i> th cell of $\Gamma_{\mu}$ .	$S_{\text{Gibbs}}(\rho) = \text{expected info gain}$ upon finding multi-particle system to be in microstate $x$ at time $t$ .

## 3. S<sub>Shan</sub> as Maximum Amount of Message Compression

- Let  $X = \{x_1, ..., x_\ell\}$  be a set of letters from which we construct messages.
- Suppose the messages have *N* letters a piece.
- Let  $\{p_1, ..., p_\ell\}$  be a probability distribution over *X*.



• <u>Thus</u>:

<u>So</u>:

$$\log_2 \left( \begin{array}{c} \text{The number of distinct} \\ \text{typical messages} \end{array} \right) = \log_2 \left( \frac{N!}{(p_1 N)! \cdots (p_\ell N)!} \right)$$

Let's simplify the RHS...

$$\log_{2} \left( \frac{N!}{(p_{1}N)! \cdots (p_{\ell}N)!} \right) = \log_{2}(N!) - \{ \log_{2}(p_{1}N)! + \cdots + \log_{2}(p_{\ell}N)! \}$$

$$Stirling's approx: \log_{2}n! \approx n\log_{2}n - n$$

$$\approx (N\log_{2}N - N) - \{ (p_{1}N\log_{2}p_{1}N - p_{1}N) + \cdots + (p_{\ell}N\log_{2}p_{\ell}N - p_{\ell}N) \}$$

$$= N\{\log_{2}N - 1 - p_{1}\log_{2}p_{1} - p_{1}\log_{2}N + p_{1} - \cdots - p_{\ell}\log_{2}p_{\ell} - p_{\ell}\log_{2}N + p_{\ell} \}$$

$$= -N\sum_{i}p_{i}\log_{2}p_{i}$$

$$= NS_{\text{Shan}}(X)$$

• Thus:  

$$\log_2 \left( \begin{array}{c} The number of distinct \\ typical messages \end{array} \right) = NS_{\text{Shan}}(X)$$
• So:  
(The number of distinct \\ typical messages ) = 2^{NS\_{\text{Shan}}(X)} \qquad \log\_2 x = y \implies x = 2^y

- <u>So</u>: There are only  $2^{NS_{\text{Shan}}(X)}$  typical messages with N letters.
- This means, at the message level, we can encode them using only  $NS_{Shan}(X)$  bits.



• <u>Now</u>: At the letter level, how many bits are needed to encode a message of *N* letters drawn from an ℓ-letter alphabet?

<i><u>First</u></i> : How many bits are needed to encode each letter in an $\ell$ -letter alphabet?			
1	$\ell = \#letters$	<u>x = #b</u>	its per letter
ł	2 letters	1 bit:	0, 1
	4 letters	2 bits:	00, 01, 10, 11
	8 letters	3 bits:	000, 001, 010, 011, 100, 101, 110, 111
<u>So</u> :	$\ell = 2^x$ , or $x = \log_2 \ell$		

- <u>Note</u>:  $\log_2 \ell$  bits per letter entails  $N \log_2 \ell$  bits for a sequence of N letters.
- <u>*Thus*</u>: If we know how probable each letter is, then instead of requiring  $N\log_2 \ell$  bits to encode our messages, we can get by with only  $NS_{Shan}(X)$  bits.
- <u>So</u>:

 $S_{\text{Shan}}(X)$  represents the maximum amount that typical messages drawn from a set of letters with a probability distribution defined on it can be compressed.

<u>*Ex*</u>: Let  $X = \{A, B, C, D\}$  ( $\ell = 4$ )

- <u>*Then*</u>: We need  $\log_2 4 = 2$  bits per letter.

$$For instance:$$
  
 $A = 00, B = 01, C = 10, D = 11$ 

- *So*: We need 2*N* bits to encode a message with *N* letters.
- <u>Now</u>: Suppose the probabilities for each letter to occur in a typical N-letter message are the following:

$$p_A = 1/2$$
,  $p_B = 1/4$ ,  $p_C = p_D = 1/8$ 

- *Then*: The minimum number of bits needed to encode all possible *N*-letter messages is:

$$NS_{\text{Shan}}(X) = -N\left(\frac{1}{2}\log_2\frac{1}{2} + \frac{1}{4}\log_2\frac{1}{4} + \frac{1}{8}\log_2\frac{1}{8} + \frac{1}{8}\log_2\frac{1}{8}\right) = 1.75N$$

- <u>*Thus*</u>: If we know how probable each letter is, instead of requiring 2*N* bits to encode all possible messages, we can get by with only 1.75*N*.
- <u>Note</u>: If all letters are equally likely (the equilibrium distribution), then  $p_A = p_B = p_C = p_D = 1/4$ .
- <u>And</u>:  $NS_{\text{Shan}}(X) = -N(\frac{1}{4}\log_2\frac{1}{4} + \frac{1}{4}\log_2\frac{1}{4} + \frac{1}{4}\log_2\frac{1}{4} + \frac{1}{4}\log_2\frac{1}{4}) = 2N$

#### <u>Message Compression Interpretation Comparison</u>

Shannon	Boltzmann
X = set of letters	$\Gamma_{\mu}$ = single-particle phase space
$\{x_1,, x_\ell\} = $ letters	$\{x_1,, x_\ell\}$ = single-particle microstates
<i>N</i> -letter message	<i>N</i> -particle microstate
N = # of letters in message	<i>N</i> = # single-particle microstates in a multi-particle microstate
$\{p_1,, p_\ell\} = \text{probability}$ distribution over letters	$\{p_1,, p_\ell\}$ = probabilty distribution over single-particle microstates
$p_i$ = probability that letter $x_i$ occurs in a message	$p_i$ = prob that single-particle microstate $x_i$ occurs in an <i>N</i> -particle microstate
$Np_i = #$ of occurrences of letter $x_i$ in typical message	$Np_i = #$ occurrences of single-particle microstate $x_i$ in typical <i>N</i> -particle microstate
$S_{\rm Shan}(X) = -\sum_i p_i \log_2 p_i$	$S_{\text{Boltz}}(\Gamma_{M_D}) = -Nk\sum_i p_i \ln p_i$
$NS_{\text{Shan}} = \text{minimum number of}$ base 2 numerals ("bits") needed to encode a message composed of <i>N</i> letters drawn from set { $x_1,, x_\ell$ }.	$S_{\text{Boltz}} \sim NS_{\text{Shan}} = \text{minimum number of}$ base <i>e</i> numerals (" <i>e</i> -bits?") needed to encode a multi-particle microstate composed of <i>N</i> single-particle microstates drawn from { $x_1,, x_\ell$ }.

Gibbs?

#### Interpretive Issues:

(1) How should the probabilities  $p_i$  in  $S_{\text{Shan}}(X) = -\sum_i p_i \log_2 p_i$  be interpreted?

- <u>Emphasis is on uncertainty</u>: The information content of the value  $x_i$  of a random variable X is a function of how uncertain it is, with respect to the receiver.
  - <u>So</u>: Perhaps the probabilities are *epistemic*.
  - In particular: p<sub>i</sub> is a measure of the receiver's degree of belief in the accuracy of the value x<sub>i</sub>.
- *But*: The probabilities are set by the nature of the source.
  - If the source is not probabilistic, then  $p_i$  can be interpreted epistemically.
  - If the source is inherently probabilistic, then  $p_i$  can be interpreted as the *ontic* probability that the source produces the value  $x_i$ .

(2) How is  $S_{\text{Shan}}$  related to other notions of entropy?

Thermodynamic	$S_{\rm TD}(\sigma_2) = \int_{\sigma_1}^{\sigma_2} \frac{\delta Q_R}{T} + S_0$
Boltzmann:	$S_{\text{Boltz}}(\Gamma_M) = k \ln  \Gamma_M $ = $-k \sum_{i=1}^{\ell} n_i \ln n_i + \text{const.}$
	$= -Nk \sum_{i=1}^{\ell} p_i \ln p_i + \text{const.}$
	$S_{\text{Boltz}}(\Gamma_M) = -Nk \int_{\Gamma_\mu} \rho_\mu(x_\mu) \ln \rho_\mu(x_\mu) dx_\mu$
Gibbs:	$S_{\text{Gibbs}}(\rho) = -k \int_{\Gamma} \rho(x) \ln \rho(x) dx$
Shannon:	$S_{\text{Shan}}(X) = -\sum_{i=1}^{\ell} p_i \log_2 p_i$
	$S_{\text{Shan}}(X) = -\int_X \rho(x)\log_2\rho(x)dx$

Can statistical mechanics be given an information-theoretic foundation? Can the 2nd Law be given an information-theoretic foundation?