

Entropy, Compression & Entanglement concentration

Shannon '48:

1. Noiseless coding theorem:

How much can a message be compressed (asymptotically)?
What is its information content?

2. Noisy channel coding theorem:

At what rate can we communicate reliably over a noisy channel?

He solved both problems and today the applications are ubiquitous. Quantumly, only the first problem is fully solved, so we will focus there.

1. Shannon entropy H :

Definition: For a probability distribution $\vec{p} \in \mathbb{R}^k$, the entropy of \vec{p} is defined as

$$H(\vec{p}) = - \sum_{i=1}^k p_i \log p_i$$

Examples:

$$- H(1, 0, 0, \dots, 0) = 0$$

$$- H\left(\frac{1}{k}, \frac{1}{k}, \dots, \frac{1}{k}\right) = \log_2 k$$

$$- \text{in general } 0 \leq H(\vec{p}) \leq \log_2 k$$

$$- H(\vec{p} \otimes \vec{q}) = H(\vec{p}) + H(\vec{q})$$

$$- \text{in general, for a joint distribution } p_{AB} \in \mathbb{R}^k \otimes \mathbb{R}^l, \\ \max\{H(p_A), H(p_B)\} \leq H(p_{AB}) \leq H(p_A) + H(p_B)$$

2. Shannon entropy and classical data compression

Stirling's approximation:

$$\lim_{n \rightarrow \infty} \frac{n!}{\sqrt{2\pi n} \left(\frac{n}{e}\right)^n} = 1$$

$$\Rightarrow \ln(n!) = \ln\left(\frac{n}{e}\right)^n + O(\log n) \\ = n \ln n - n + O(\log n)$$

Corollary: For $p \in [0, 1]$,

$$\log \binom{n}{pn} = \log \frac{n!}{(pn)! [(1-p)n]!}$$

$$\approx n \log n - (\log e) n - [pn \log(pn) - (\log e) pn + (1-p)n \log(1-p)n - \log e(1-p)n] \\ = n \cdot H(p).$$

More generally,

Corollary: For \vec{p} a distribution on k elements,

$$\log \frac{n!}{\prod_{i=1}^k (p_i n)!} \approx n \cdot H(\vec{p}).$$

By the law of large numbers, if we draw n independent samples from \vec{p} (ie, one sample from $\vec{p}^{\otimes n}$), with probability exponentially close to 1, every outcome i will be observed about $p_i n$ times.

times.

The number of strings in $\{1, 2, \dots, k\}^n$ with $p_i n$ i 's is

$$\frac{n!}{\prod_{i=1}^k (p_i n)!}$$

These "typical" strings can be compressed using $n \cdot H(\vec{p})$ bits of information.

More precisely,

(for all $\epsilon > 0$ and n large enough)

Noiseless coding theorem: There exists a set $C \subseteq \{1, \dots, k\}^n$ of size $|C| \leq 2^{n(H(\vec{p}) + \epsilon)}$ such that

$$\mathbb{P}[X \in C] \geq 1 - \epsilon$$

when sample X is drawn from the distribution $\vec{p}^{\otimes n}$.

Proof:

$$\mathbb{P}[x_1, x_2, \dots, x_n] = p(x_1) p(x_2) \dots p(x_n)$$

$$\therefore \log \mathbb{P}[x_1, \dots, x_n] = \sum_{i=1}^n \log p(x_i)$$

$$\therefore \log \mathbb{P}[X] = \sum_{i=1}^n \log p(X_i) \quad \begin{array}{l} \text{sum of independent r.v.s} \\ \text{with finite mean \& varrrance} \end{array}$$

By the Central Limit Theorem, this sum concentrates to its expectation, $n \cdot \sum_{j=1}^k p_j \log p_j = -n H(\vec{p})$. We may say

$$\mathbb{P}\left[\left| \frac{1}{n} \log \mathbb{P}[X] - H(\vec{p}) \right| \leq \delta\right] \geq 1 - \epsilon$$

\therefore Let $C = \{x : \mathbb{P}[x] \in [2^{-n(H+\delta)}, 2^{-n(H-\delta)}]\}$. Then

$$\mathbb{P}[X \in C] \geq 1 - \epsilon$$

and, necessarily, $|C| \leq 2^{n(H+\delta)}$ since the total probability is at most one. \square

Moreover, this construction is asymptotically optimal, since with fewer than $n(H-\delta)$ bits we could not cover all typical sequences, if $\delta' > \delta$,

$$\mathbb{P}[\text{successful decoding}] \leq 2^{n(H-\delta')} 2^{-n(H-\delta)} + \epsilon = 2^{-n(\delta'-\delta)} + \epsilon$$

3. Von Neumann entropy S

Definition: For a quantum state ρ , that can be diagonalized as $\sum_i p_i |\psi_i\rangle\langle\psi_i|$ with $\langle\psi_i|\psi_j\rangle = \delta_{ij}$, the Von Neumann entropy of ρ is

$$S(\rho) = H(\vec{p}) = -\text{Tr}(\rho \log \rho).$$

Properties:

a. $S(|\psi\rangle\langle\psi|) = 0$

b. $S(U\rho U^\dagger) = S(\rho)$

c. $S(\rho) \leq \log D$, with equality only for the maximally mixed state $\rho = \frac{1}{D} \mathbb{1}$.

d. Concavity: for $\lambda \in [0, 1]$,
 $S(\lambda\rho_1 + (1-\lambda)\rho_2) \geq \lambda S(\rho_1) + (1-\lambda)S(\rho_2)$.

Interpretation: If we are more ignorant of how a state was prepared, then the entropy is higher.

e. Subadditive: $S(\rho_{AB}) \leq S(\rho_A) + S(\rho_B)$.

f. Strongly subadditive:

$$S(\rho_{ABC}) + S(\rho_B) \leq S(\rho_{AB}) + S(\rho_{BC})$$

(difficult to prove, but extremely useful)

* g. Δ inequality:

$$S(\rho_{AB}) \geq |S(\rho_A) - S(\rho_B)|$$

Unlike for the Shannon entropy, it is not true that $S(\rho_{AB}) \geq \max\{S(\rho_A), S(\rho_B)\}$.

Example: If $\rho_{AB} = |\psi\rangle\langle\psi|_{AB}$ is a pure state, then $S(\rho_{AB}) = 0$, but

$$S(\rho_A) = S(\rho_B) \text{ may be } > 0.$$

↑
by Schmidt decomposition

Unlike classically, discarding a subsystem can increase your uncertainty about the state.

More operational properties:

h. If Y is the result of any full measurement of ρ , i.e. a classical random variable with $P(Y=i) = p_{ii}$, then $H(Y) \geq S(\rho)$, with equality iff the measurement basis diagonalizes ρ .

i. Preparing the ensemble $\{p_x, |\psi_x\rangle\}$, $H(X) \geq S\left(\sum_x p_x |\psi_x\rangle\langle\psi_x|\right)$