# Entropy & Information

## Shannon Entropy — Recap!

Suppose we learn a random variable $X$

$$H(X) = -\sum_x p_x \log p_x \equiv \text{"uncertainty about } X \text{ before learning it"}$$

$\equiv$ "Amount of information we gain on learning $X$"

Note $\lim_{x \to} x \log x = 0$

$\log \equiv \log_2$

Suppose a source is producing data in the form of random variables $X_1, X_2, X_3 \ldots$

Suppose each random variable can take a character $x_u$ with probability $p_u$.

What's the minimal physical resources required to store the data produced by the source?

Answ: $n$ symbol string can be compressed to $n\,H(X)$ symbols

Shannon's noiseless coding theorem

eg. Suppose a source of information produces $1, 2, 3 \sim 4$
with probabilities $\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8}$

A naive binary encoding $1 = 00 \quad 2 = 01 \quad 3 = 10 \quad 4 = 11$

On average the length of a string with this encoding is

$$2 \times \frac{1}{2} + 2 \times \frac{1}{4} + 2 \times \frac{1}{8} + 2 \times \frac{1}{8} = 2$$

Then we can use the bias to reduce the amount of symbols required to store data from that source by using less characters to store commonly obtained symbols & more to store less likely ones.

eg.     1 = 0     2 = 10     3 = 110     4 = 111

On average the length of a string with this encoding is

$$1 \times \frac{1}{2} + 2 \times \frac{1}{4} + 3 \times \frac{1}{8} + 3 \times \frac{1}{8} = \frac{7}{4} \leq 2$$

New code is more efficient!

Sketch of general proof:

Binary case first —     $X = 1$     with prob $p$
          $= 0$               $1-p$

Consider a data string of length $n$.

In the limit of large $n$, a typical bit string will contain about $n(1-p)$ 0s and $np$ 1s.

There are $^nC_{np}$ such typical strings

$$\log \left( ^nC_{np} \right) = \log \left( \frac{n!}{(np)! \; n(1-p)!} \right)$$

$$= \log n! - \log(np!) - \log(n(1-p)!)$$

Use Stirling approximation  $\log(n!) \simeq n\log n - n$  (in limit of large $n$)

$$\log \left( ^nC_{np} \right) \simeq n\log n - n - (np \log np - np + n(1-p)\log(n(1-p)) - n(1-p))$$

$$= -np \log p - n(1-p) \log(1-p) = n H(p)$$

$$\Rightarrow \quad \text{no. of typical strings} \simeq 2^{n H(p)} \quad \text{binary entropy}$$

Compression strategy — assign a positive integer to each of the possible typical bit strings.

There are $2^{n H(p)}$ such strings

so $2^{n H(p)}$ letters are required

& each letter can be encoded using $n H(p)$ bits.

Note — the completely uniform distribution cannot be compressed.

ie. $H(\tfrac{1}{2}) = -\tfrac{1}{2} \log(\tfrac{1}{2}) - \tfrac{1}{2} \log(\tfrac{1}{2}) = -\log(\tfrac{1}{2}) = \log 2 = 1$

i.e. $n$ bits are 'encoded' in $n$ bits

Generalisation beyond binary case.

If letter $k$ occurs with probability $p_k$ in a string of length $n$ each $k$ will typically occur $n p_k$ times

There are $\dfrac{n!}{\prod_k (n p_k)!}$ such typical strings

& $\dfrac{n!}{\prod_k (n p_k)!} \simeq 2^{n H(X)} \quad \Rightarrow \quad n H(X)$ binary encoding possible

$\Rightarrow$ Operational interpretation of Shannon entropy!

# Conditional Entropy & Mutual Information

Consider 2 random variables $X$ & $Y$

– How is the information content of $X$ related to $Y$?

Conditional entropy & Mutual information provide answers

But first: **Joint Entropy**

$$H(X,Y) = \sum_{x,y} p(x,y) \log(p(x,y))$$

This is the total uncertainty about $X$ & $Y$

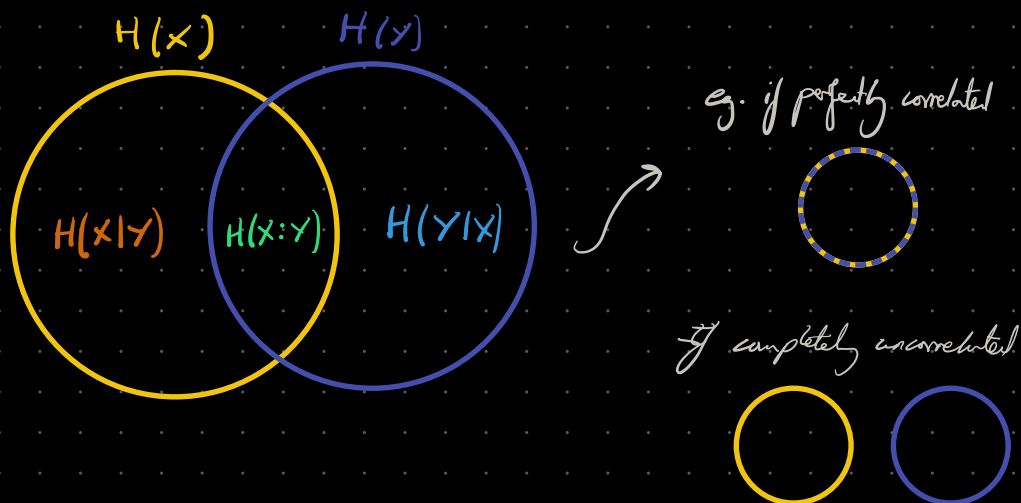Suppose we know the value of $Y$, so we have gained $H(Y)$ bits of information.

The **conditional entropy** of $X$ on knowing $Y$ is the remaining uncertainty in $X$ on knowing $Y$.

$$H(X|Y) = H(X,Y) - H(Y)$$

**Mutual information** measures the amount of information $X$ and $Y$ have in common – ie. measures their correlations

$$H(X:Y) = H(X) + H(Y) - H(X,Y)$$

The following Venn diagram is a super useful tool to get a sense of their properties.

$H(x)$　　$H(y)$



eg. if perfectly correlated

If completely uncorrelated

Can read off some properties straight from the venn:

- $0 \leq H(X|Y) \leq H(X)$

- $H(X|Y) \neq H(Y|X)$

- $0 \leq H(X:Y) \leq \min\{H(X), H(Y)\}$

Drawing Venn is helpful for providing an intuition but is not the full story — always prove inequalities also independently.
(Problem sheet for this week will provide many)

**Relative Entropy** — a measure of the closeness of 2 distributions. Useful for proving stuff.

$$H(p(x) \| q(x)) = \sum_x p(x) \log\left(\frac{p(x)}{q(x)}\right) \longrightarrow = 0 \text{ if } p(x) = q(x) \; \forall x$$

$$= -H(X) - \sum_x p(x) \log q(x)$$

- $H(p(x) \| q(x)) \geq 0$　　( to prove this use $-\log(x) \geq \frac{1-x}{\ln 2}$
$$\Rightarrow H(p(x) \| q(x)) \geq \frac{1}{\ln 2} \sum_x p(x) \left(1 - \frac{q(x)}{p(x)}\right) )$$
$$\underbrace{}_{= 1 - 1 = 0}$$

- $H(p(x) \| \frac{1}{d}) = -H(X) - \sum_x p(x) \log 1/d$
$$= \log(d) - H(X)$$

Shannon entropy $\equiv$ Relative entropy to max. uncertain distribution

# Von Neumann Entropy

$$S(\rho) = -\text{Tr}(\rho \log \rho) = -\sum_i \lambda_i \log \lambda_i$$

eigs of $\rho$

$\nearrow$

$$= H(\{\lambda_i\})$$

Similarly to classical case $S(\rho)$ quantifies the compressibility of quantum information

$\rho^{\otimes n}$ can be compressed to $\sigma$ that lives on a Hilbert space $H_C$

with $\dim(H_C) = 2^{n S(\rho)}$

~~Intuition~~ roughly the same $\rho$ can look at subspace corresponding to "typical" sequences of eigenvalues.

See Preskill's notes for a proof.

# Important Properties

1) Pure states have zero entropy

$$\rho = |\psi\rangle\langle\psi| \qquad \lambda = 1 \qquad S(\rho) = \log(1) = 0$$

2) Invariance: $\quad S(U\rho U^\dagger) = S(\rho)$ (eigenvalues are left unchanged)

3) Maximum: $\quad \text{Max } S(\rho) = S(\mathbb{1}/d) = \log(d)$

4) Entropy of measurement:

Say you measure $M = \{m_j |m_j\rangle\langle m_j|\}$

$p(m_j) = \langle m_j| \rho |m_j\rangle$

$Y = \{m_j, p(m_j)\}$

$\Rightarrow \quad H(Y) \geq S(\rho)$

Equivalent to the statement that replacing $\rho$ in any basis with its decohered variant increases entropy.

ie. Killing off coherences increases entropy.

5) Additivity: $S(\rho_A \otimes \rho_B) = S(\rho_A) + S(\rho_B)$

"eigenvalues multiply — take log — entropies add"

6) Triangle Inequality

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

Use Klein's Inequality: $S(\rho) \leq -Tr(\rho \log \sigma)$

Let $\rho = \rho_{AB}$ & $\sigma = \rho_A \otimes \rho_B \rightsquigarrow$ $S(\rho) \leq -Tr(\rho^{AB}(\log(\rho^A) + \log(\rho^B)))$

$= -Tr(\rho^A \log \rho^A) - Tr(\rho^B \log \rho^B)$

$= S(\rho_A) + S(\rho_B)$

7) Concavity: $S(\sum_i p_i \rho_i) \geq \sum_i p_i S(\rho_i)$

"extra randomness only increases uncertainty"

$\rho_i = \sum_j \lambda_j^i |\lambda_j^i\rangle\langle\lambda_j^i|$

Let $\rho_{AB} = \sum_i p_i \rho_i \otimes |i\rangle\langle i|$

$\rho_A = \sum_i p_i \rho_i$ $\qquad \rho_B = \sum_i p_i |i\rangle\langle i|$

$S(\sum_{ij} p_i \lambda_j^i |\lambda_j^i\rangle\langle\lambda_j^i|)$

$= -\sum_{ij} p_i \lambda_j^i \log p_i \lambda_j^i$

$= -\sum_{ij} p_i \lambda_j^i \log p_i + \sum_{ij} p_i \lambda_j^i \log \lambda_j^i$

$H(\{p\}) + \sum_i p_i S(\rho_i)$

$S(\rho_{AB}) \leq S(\rho_A) + S(\rho_B)$ $\Rightarrow$ $H(\{p\}) + \sum_i p_i S(\rho_i) \leq S(\sum_i p_i \rho_i)$

$\underset{S(\sum_i p_i \rho_i)}{\overset{\downarrow}{\rule{0pt}{1em}}} \underset{H(\{p\})}{\overset{\hookrightarrow}{\rule{0pt}{1em}}}$ $\qquad + H(\{p\})$

$\checkmark$

Analagously to the classical case we can define :

Joint entropy $\qquad S(\rho_{AB}) = -\text{Tr}(\rho_{AB}\log(\rho_{AB}))$

Conditional entropy $\qquad S(\rho_A|\rho_B) = S(\rho_{AB}) - S(\rho_B)$

Mutual information $\qquad S(\rho_A : \rho_B) = S(\rho_A) + S(\rho_B) - S(\rho_{AB})$

Note that the Venn diagram breaks down in this case

eg. Conditional entropy can be negative

$\qquad$ Say $\quad \rho_{AB} = |\phi^+\rangle\langle\phi^+| \qquad \rho_A = \rho_B = \frac{1}{2}$
$\qquad\qquad S(\rho_{AB}) = 0 \qquad\qquad S(\rho_A) = -\frac{1}{2}\log\frac{1}{2} - \frac{1}{2}\log\frac{1}{2} = \log(2) = 1$

$\qquad\qquad S(\rho_A | \rho_B) = -1 ! \qquad$ "Uncertainty in joint state is less
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ than the reduced states"

## Relative Entropy

- Not a distance measure [3]
- But can be used to measure the similarity between two quantum states [1] [2]

$$S(\rho \| \sigma) := Tr\left( \rho \log \rho - \rho \log \sigma \right)$$

$$\overset{\Sigma_i \lambda_i |\lambda_i\rangle\langle\lambda_i|}{\quad} \qquad \overset{\Sigma_j \mu_j |\mu_j\rangle\langle\mu_j|}{\quad}$$

$$= \Sigma_i \lambda_i \left( \log \lambda_i - \sum_j |\langle \mu_j | \lambda_i \rangle|^2 \log(\mu_j) \right)$$

Reduces to the classical relative entropy if diagonal in same basis but more generally depends on the overlap between their eigenbasis

Properties:

1) Positivity: $S(\rho \| \sigma) \geq 0$

2) Faithful: $S(\rho \| \sigma) = 0$ iff $\rho = \sigma$

3) Aysymetric: $S(\rho \| \sigma) \neq S(\sigma \| \rho)$

4) Unitarily invariant: $S(U \rho U^\dagger \| U \sigma U^\dagger) = S(\rho \| \sigma)$

clear from

# Data processing Inequality

= holds for Unitary evolutions

$$S(\, \mathcal{E}(\rho) \,\|\, \mathcal{E}(\sigma)\,) \;\leq\; S(\rho \| \sigma) \qquad \forall \, \mathcal{E}$$

" There is no channel you can apply that will make $\rho$ & $\sigma$ more distinguishable "

Same also holds for 1 norm

$$\|\, \mathcal{E}(\rho) - \mathcal{E}(\sigma)\,\|_1 \;\leq\; \|\rho - \sigma\|_1 \qquad \forall \, \mathcal{E}$$

(But it doesn't hold for 2-norm)

# Mixed States Fidelity

$$F(\rho, \sigma) = Tr\left(\sqrt{\rho^{\frac{1}{2}} \sigma \rho^{\frac{1}{2}}}\right)$$

Case 1 :  $\rho$  &  $\sigma$  commute  $\Rightarrow$   $\rho = \sum_i r_i |i\rangle\langle i|$   $\sigma = \sum_i s_i |i\rangle\langle i|$

$$F(\rho, \sigma) = Tr\left(\sqrt{\sum_i r_i s_i |i\rangle\langle i|}\right)$$

$$= Tr\left(\sum_i \sqrt{r_i s_i} |i\rangle\langle i|\right)$$

$$= \sum_i \sqrt{r_i s_i}$$

$$= F(r, s) \checkmark \quad \text{Classical fidelity}$$

Case 2 :   $\sigma$ ,  $\rho = |\psi\rangle\langle\psi|$     $(|\psi\rangle\langle\psi|)^2 = \rho \Rightarrow \rho^{\frac{1}{2}} = |\psi\rangle\langle\psi|$

$$F(\rho, \sigma) = Tr\left(\sqrt{|\psi\rangle\langle\psi| \sigma |\psi\rangle\langle\psi|}\right)$$

$$= Tr\left(\sqrt{\langle\psi|\sigma|\psi\rangle} \sqrt{|\psi\rangle\langle\psi|}\right)$$

$$|\psi\rangle\langle\psi|$$

$$= \sqrt{\langle\psi|\sigma|\psi\rangle} \quad \longleftarrow \quad \begin{array}{l}\text{Fidelity between pure and}\\\text{mixed state is equal to}\\\text{the overlap}\end{array}$$

Case 2b.    $\sigma = |\phi\rangle\langle\phi|$      $F(\rho, \sigma) = |\langle\psi|\phi\rangle|$

Note the lack of mod
square here → this is
a matter of convention
I'm following N&C here.

General case?   Operational interpretation provided by   Uhlmann's Theorem.

**Uhlmann's Theorem:**

$$F(\rho, \sigma) = \max_{|\psi\rangle, |\varphi\rangle} |\langle \psi | \varphi \rangle|$$

max over all possible purifications of $\psi$ & $\varphi$.

where $\rho_S = \text{Tr}_R (|\psi\rangle\langle\psi|_{RS})$ & $\sigma_S = \text{Tr}_R (|\varphi\rangle\langle\varphi|_{RS})$

proof — exercise sheet this week.

Data processing inequality also holds here

$$F(\mathcal{E}(\rho), \mathcal{E}(\sigma)) \geq F(\rho, \sigma)$$