

Entropy:

alphabet



↳ discrete random variable: $X, A = \{x_1, \dots, x_n, \dots\}$

↳ probability function of X : $p(x) \geq 0, \sum_i p(x_i) = 1$.

* The entropy of X is a measure of the uncertainty of a random variable defined as: $H(X) = - \sum p(x_i) \log p(x_i) \geq 0$
 ↑ ↑
 in bits base 2 log

* Various possible axiomatic justifications. Shannon:
 ↳ sequence of symmetric functions $H_m(p_1, \dots, p_m)$
 such that i) $H_2(1/2, 1/2) = 1$, ii) $H(p, 1-p)$ continuous in p
 iii) $H_m(p_1, \dots, p_m) = H_{m-1}(p_1+p_2, p_3, \dots, p_m) + (p_1+p_2) H_2(\frac{p_1}{p_1+p_2}, \frac{p_2}{p_1+p_2})$

* The joint entropy of (X, Y) is $H(X, Y) = - E[\log p(X, Y)]$

* The conditional entropy of $Y|X$ is
 $H(Y|X) = \sum_x p(x) H(Y|X=x) = - \sum_x p(x) \sum_y p(y|x) \log p(y|x)$

* The chain rule
 $H(X, Y) = - \sum_x \sum_y p(x, y) \log(p(x) p(y|x)) = - \sum_x p(x) \log p(x) - \sum_x p(x) \sum_y p(y|x) \log p(y|x)$
 $= H(X) + H(Y|X)$

Generalization:
 $H(x_1, \dots, x_n) = \sum_{i=1}^n H(x_i | x_{i-1}, \dots, x_1) \leq \sum_{i=1}^n H(x_i)$
 ↳ nested conditioning not independence bound

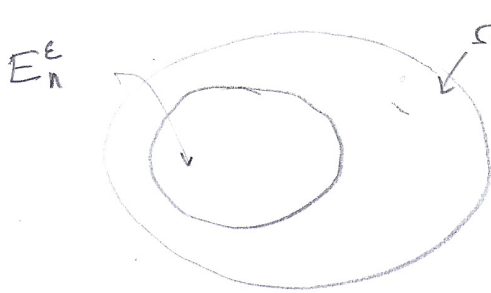
Entropy and description length

x_1, \dots, x_n i.i.d then the asymptotic equipartition property (law of large numbers) states that

$$-\frac{1}{n} \log p(x_1, \dots, x_n) \xrightarrow{n \rightarrow \infty} H(X)$$

"Almost all events are almost equally surprising"

Shannon code: typical set: $E_n^\epsilon: \{x_1, \dots, x_n \mid 2^{-n(H(X)+\epsilon)} \leq p(x) \leq 2^{-n(H(X)-\epsilon)}\}$



$$1 = \sum_x p(x) \geq \sum_{x \in E_n^\epsilon} p(x) \geq |E_n^\epsilon| 2^{-n(H(X)+\epsilon)}$$

$$\Rightarrow |E_n^\epsilon| \leq 2^{n(H(X)+\epsilon)}$$

1 bit to ensure integers
1 bit to encode whether $x \in E_n^\epsilon$

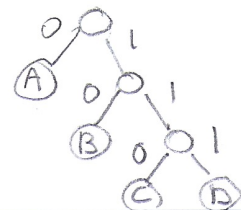
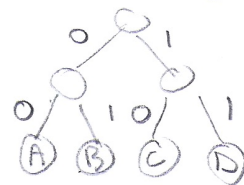
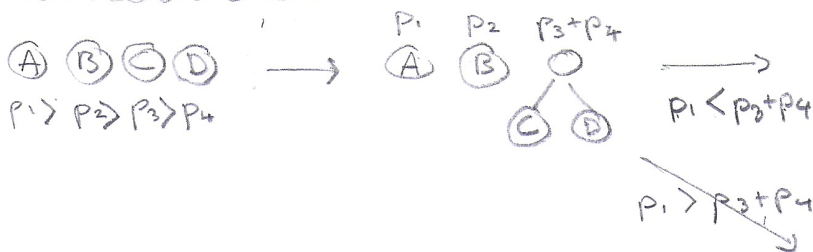
Description length of E_n^ϵ : $n(H(X)+\epsilon) + 2$ bits

of $\Omega \setminus E_n^\epsilon$: $n(\log A) + 2$ bits

Average length: $E[L(X^n)] = P[E_n^\epsilon] [n(H(X)+\epsilon) + 2] + (1 - P[E_n^\epsilon]) [n(\log A) + 2]$

n large enough $\leq n(H(X)+\epsilon) + 2 + \epsilon n(\log A) + 2$
 $= \tilde{n}(H(X)+\epsilon) + o(1)$

Huffman code: $H(X)$: optimal average code length.



$$H(X) \leq E[L] < H(X) + 1$$

E expected code length

Kullback - Leibler divergence (relative entropy)

(3)

↳ KL divergence measure how much two distributions p and q differ from one another. Informally, it quantifies the inefficiency of coding p while using an optimal code for q .

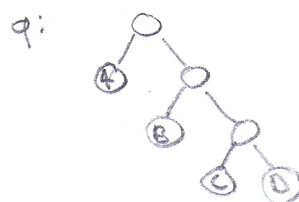
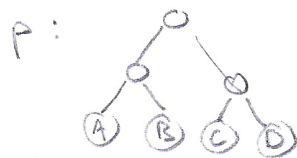
* The KL divergence between p and q is $D[p||q] = \mathbb{E}_p[\log \frac{p}{q}] \geq 0$

* $D[p||q]$ is convex in the pair (p, q) :

$$D[\lambda p_1 + (1-\lambda)p_2 || \lambda q_1 + (1-\lambda)q_2] \leq \lambda D[p_1||q_1] + (1-\lambda)D[p_2||q_2]$$

This follows from Jensen inequality + convexity of $-\log$

* Huffman trees interpretation



$$D_{KL}(p||q) = \underbrace{-\mathbb{E}_p[\log q]}_{\substack{\uparrow \text{average code length using optimal encoding for } q}} + \underbrace{H(p)}_{\substack{\uparrow \text{optimal description length of } p}}$$

* The chain rule: $D[p(x, y)||q(x, y)] = D(p(x)||q(x)) + D(p(y|x)||q(y|x))$

* "Second law of thermodynamics":

$$p(x_n, x_{n+1}) = K(x_{n+1}|x_n)p(x_n)$$

$$q(x_n, x_{n+1}) = K(x_{n+1}|x_n)q(x_n)$$

$$D[p(x_n, x_{n+1})||q(x_n, x_{n+1})]$$

≥ 0

\uparrow Markov kernel

$$= D[p(x_n)||q(x_n)] + D[K(x_{n+1}|x_n)||K(x_{n+1}|x_n)]$$

$$= D[p(x_{n+1})||q(x_{n+1})] + D[p(x_n|x_{n+1})||q(x_n|x_{n+1})] \geq 0$$

$$\Rightarrow D[p(x_{n+1})||q(x_{n+1})] \leq D[p(x_n)||q(x_n)]$$

\Rightarrow KL divergence to stationary distribution decreases with time

Mutual information

(4)

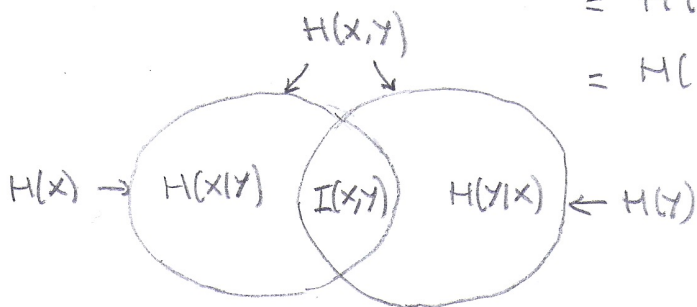
The mutual information between X and Y , denoted $I(X, Y)$ is the KL divergence between the joint distribution and the product distribution:

$$I(X, Y) = D[p(x, y) \parallel \underset{\substack{\uparrow \\ \text{marginals}}}{p(x)p(y)}] = \sum_x \sum_y p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \geq 0$$

$$I(X, Y) = 0 \quad (\Leftrightarrow) \quad X, Y \text{ independent}$$

↳ conditioning reduces entropy

* Alternatively: $I(X, Y) = H(X) - H(X|Y)$ ← reduction of uncertainty of X (resp Y) due to knowledge of Y (resp X).
 $= H(Y) - H(Y|X)$
 $= H(X) + H(Y) - H(X, Y)$



* Conditional mutual information

$$I(X; Y|Z) = H(X|Z) - H(X|Y, Z)$$

$$= \mathbb{E}_{p(x, y, z)} \left[\log \frac{p(x, y|z)}{p(x|z)p(y|z)} \right]$$

$$I(X; Y|Z) = 0 \quad X|Z, Y|Z \text{ independent}$$

* The chain rule: $I(X_1, \dots, X_n; Y) = \sum_{i=1}^n I(X_i, Y | X_{1:i-1})$

* $I(X, Y)$ is concave in $p(x)$ and convex in $p(y|x)$

↳ concavity: $I(X, Y) = H(Y) - \sum_x p(x) H(Y|X=x)$
 ↑ fixed

H concave in $p(y)$ ← linear in $p(x)$

↳ convexity: consider $x \begin{cases} \xrightarrow{p_1(y|x)} Y_1 \\ \xrightarrow{p_2(y|x)} Y_2 \end{cases}$ for fixed $p(x)$

consider $Y^{(\lambda)} = Y_\theta$ where $\theta = \begin{cases} 1 & \text{with probability } \lambda \\ 2 & \text{otherwise} \end{cases}$
 independent \uparrow

$x \rightarrow Y^{(\lambda)}$ with conditional probability
 $p_\lambda(y|x) = \lambda p_1(y|x) + (1-\lambda) p_2(y|x)$

we can check that: $p_\lambda(x,y) = \lambda p_1(x,y) + (1-\lambda) p_2(x,y)$
 $p(x)p_\lambda(y) = \lambda p(x)p_1(y) + (1-\lambda) p(x)p_2(y)$

$I(Y^{(\lambda)}, X) = D \left[\lambda p_1(x,y) + (1-\lambda) p_2(x,y) \parallel \lambda p(x)p_1(y) + (1-\lambda) p(x)p_2(y) \right]$
 \uparrow with mixture kernel p_λ \uparrow Kullback-Leibler convex \checkmark

* $x \rightarrow y$ encoding channel
 $p(x) \quad \uparrow \quad p(y|x)$

$\min_{p(y|x)} I(x,y) = 0 \leftarrow$ uninteresting without constraints

$\max_{p(x)} I(x,y) = C_{p(y|x)} \leftarrow$ channel capacity

Interest for neuroscience: optimizing over $p(y|x)$
 not necessarily well-posed!

Data processing inequality

(6)

Markov chain: $X \rightarrow Y \rightarrow Z$; $p(x, y, z) = p(x) p(y|x) p(z|y)$

Remark: • $X|Y$ and $Z|Y$ are independent

$$p(x, z|y) = \frac{p(x, y) p(z|y)}{p(y)} = p(x|y) p(z|y)$$

• $X \rightarrow Y \rightarrow Z \Rightarrow Z \rightarrow Y \rightarrow X$

No processing, deterministic or random, can increase the information that Y contains about X :

$$X \rightarrow Y \rightarrow Z \Rightarrow I(X, Y) \geq I(X, Z)$$

Proof:
$$\begin{aligned} I(X; Y, Z) &= I(X; Y) + I(X; Z|Y) \\ &= I(X; Z) + I(X; Y|Z) \end{aligned}$$

conditional independence $\Rightarrow I(X; Z|Y) = 0$

$$I(X; Y|Z) \geq 0 \Rightarrow I(X; Y) \geq I(X, Z) \quad \checkmark$$

Moreover: $I(X; Y|Z) \leq I(X; Y)$ if $X \rightarrow Y \rightarrow Z$

Observation of "downstream" variable reduces the dependence between "upstream" variables.