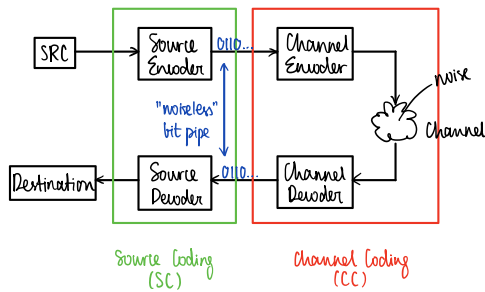


Information Theory

Big Picture



Separation Theorem: SC & CC can be done separately w/o loss of optimality from end-to-end.

SC: Impossible to compress a source X^N below its entropy $H(X^N) = NH(X)$ bits.
Possible to compress X^N to $N(H(X) + \epsilon)$ bits $\forall \epsilon > 0$ as $N \rightarrow \infty$.

CC: Impossible to transmit reliably at rate $R > C$, where C is channel capacity.
Any rate $R < C$ is achievable w/ high reliability (i.e. $\text{P}_{\text{error}}(N) \rightarrow 0$ as $N \rightarrow \infty$).

$X \sim B(p)$

If $p = \frac{1}{2}$, how much info does a single toss provide? 1 bit

What if $p = 0.11$? $\frac{1}{2}$ bit

Suppose you have a seq of indep coin tosses by Alice & Bob using fewest # of bits.
How do you do it?

Entropy

$X \sim$ discrete RV, $x \in X$

$$H(X) = E[-\log p(x)] = E[\log \frac{1}{p(x)}] = \sum_{x \in X} p(x) \log \frac{1}{p(x)} \quad (\log \text{ base } 2)$$

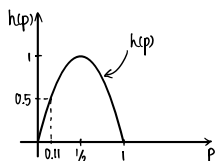
Remark: $H(X)$ measures the avg uncertainty / surprise level / information content.

Intuition: The higher $p(x)$, the lower the surprise / information you gain when you see it.

Ex: $X \sim B(p)$

$$H(X) = \sum_{x=0}^1 p(x) \log \frac{1}{p(x)} = p \log \frac{1}{p} + (1-p) \log \frac{1}{1-p}$$

$h(p)$
Binary entropy function



Some properties of entropy:

If $X \in \{1, 2, \dots, D\}$, $H(X)$ is max when it is uniformly distributed (i.e. $p(1) = p(2) = \dots = p(D)$)
 $\Rightarrow H_{\text{max}}(X) = \log_2 D$ bits $\Rightarrow H(X) \leq \log_2 D$ bits

Joint Entropy

$$H(X, Y) = \sum_x \sum_y p_{xy}(x, y) \log_2 \frac{1}{p_{xy}(x, y)}$$

$$H(X, Y) = \mathbb{E} \left[\log \frac{1}{p_{XY}(x, y)} \right]$$

Note: If X, Y indep, $H(X, Y) = H(X) + H(Y)$.

$$\begin{aligned} \text{Proof: } H(X, Y) &= -\mathbb{E}[\log(p(X, Y))] = -\mathbb{E}[\log(p(X), p(Y))] = -\mathbb{E}[\log p(X)] - \mathbb{E}[\log p(Y)] \\ &= H(X) + H(Y) \end{aligned}$$

If you observe 2 indep RVs, the info content should be additive in the info-content of each RV.

Conditional Entropy

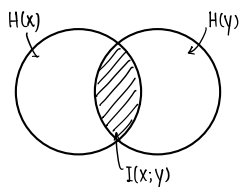
$$H(Y|X) = \mathbb{E}[\log \frac{1}{p_{Y|X}(y|x)}]$$

can show: $H(Y|X) = H(X, Y) - H(X)$
 uncertainty in Y after observing X

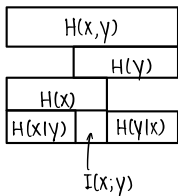
Mutual Information (MI)

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

↳ intuition: X, Y are coupled, filter out some info after observing the output



Note: If X, Y are indep, $I(X; Y) = 0$



Asymptotic Equipartition Property (AEP)

↳ Info theory version of WLLN

$$\text{Theorem: If } X_1, \dots, X_n \text{ are iid } \sim p(X), \text{ then}$$

$$\frac{-\log p(X_1, \dots, X_n)}{n} \xrightarrow{P} H(X)$$

Proof: If X, Y are indep RVs, so are $f(X)$ & $g(Y)$ for any $f(\cdot), g(\cdot)$

⇒ If X_1, \dots, X_n are indep, so are $\log p(X_1), \log p(X_2), \dots, \log p(X_n)$

$$\Rightarrow \text{WLLN: } -\frac{1}{n} \log p(X_1, \dots, X_n) = -\frac{1}{n} \sum_{i=1}^n \log p(X_i) \xrightarrow{P} -\underbrace{\mathbb{E}[\log p(X)]}_{H(X)} \text{ in prob}$$

Intuition of AEP

Ex: You flip a coin w/ bias p , n times indep

What's the prob of seeing a typical sequence?

Q/ What's the typical seq?

↳ One having np heads & $n(1-p)$ tails

Q/ What's the prob of seeing a particular typical seq?

$$\begin{aligned} P(\text{typical seq}) &= p^n (1-p)^{n(1-p)} = 2^{np \log p} 2^{nq \log q} \quad \text{where } q = 1-p \\ &= 2^{\underbrace{np \log p + nq \log q}_{-nH(p)}} \end{aligned}$$

$$P(\text{typical seq}) = 2^{-nH(p)}$$

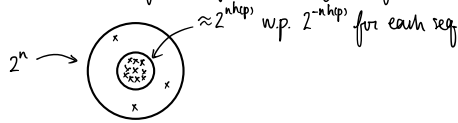
Q/ How many typical seq are there?

$$\binom{n}{np} \approx 2^{nH(p)}$$

Use Stirling's approx: $n! \approx (\frac{n}{e})^n$

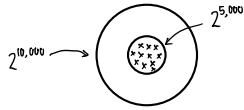
$$\frac{n!}{(np)!(nq)!} = \frac{(\frac{n}{e})^n}{(\frac{n}{e})^{np} (\frac{n}{e})^{nq}} = \frac{(\frac{n}{e})^n}{(\frac{n}{e})^{np+nq}} = \frac{(\frac{n}{e})^n}{(\frac{n}{e})^n} = 2^{n h(p)}$$

There are about $2^{n h(p)}$ typical seq. each having an equal prob of $2^{-n h(p)}$



Ex: $p = 0.11 \Rightarrow p(\text{typ seq}) = 2^{-n/2}$

$n = 10,000$



Typical Set

$A_\epsilon^{(n)}$ w.r.t. $p(x)$ is the set of sequences $(x_1, \dots, x_n) \in X^n$ s.t.

$$\Pr(2^{-n[H(x)+\epsilon]} \leq p(x_1, \dots, x_n) \leq 2^{-n[H(x)-\epsilon]}) \xrightarrow{n \rightarrow \infty} 1 \quad \forall \epsilon > 0$$

Fact: $\forall \epsilon > 0$, more than $(1-\epsilon)$ of the pmf of (x_1, \dots, x_n) lies in a typical set $A_\epsilon^{(n)}$ having at most $2^{n[H(x)+\epsilon]}$ elems.

Ex: $\epsilon = 0.001, n = 10,000, X \sim B(p=0.11)$

More than 99.9% of the prob mass lies in a set having no more than 2^{5010} elements.

Huffman Code

Prefix free code: No codeword is prefix of another.

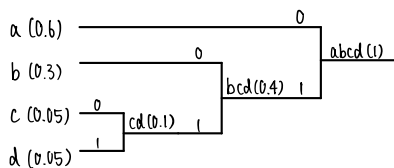
→ prevent ambiguities from occurring during decoding

→ aka "instantaneous codes" since we can decode while streaming

→ "self punctuating"

$X = \{a, b, c, d\}, P_x = \{p_a, p_b, p_c, p_d\}$

Ex: $p_a = 0.6, 0.3, 0.05, 0.05$
(a) (b) (c) (d)



Sym	Code
a	0
b	10
c	110
d	111

$\mathbb{E}[L_H] = \sum p_x l_x = 1.5$

$H(x) = \sum p_x \log \frac{1}{p_x} = 1.395$

Fact: $H(x) \leq \mathbb{E}[L_H] < H(x) + 1$

$H(x^n) \leq \mathbb{E}[L_H(x^n)] < H(x^n) + 1$

$nH(x) \leq \mathbb{E}[L_H(x^n)] < nH(x) + 1$

$H(x) \leq \frac{\mathbb{E}[L_H(x^n)]}{n} < H(x) + \frac{1}{n}$

→ Problem: Complexity