

Information Theory and Coding (Lec-02)



Eftekhari Hossain
Lecturer
Dept. of ETE, CUET

Outline

- ➔ Average Self Information
- ➔ Entropy
- ➔ Conditional Entropy
- ➔ Joint Entropy

Average Self Information

- Consider a discrete random variable X with possible outcomes x_i , $i = 1, 2, 3, \dots, n$
- The average self information of the event $X = x_i$ is defined as

$$\begin{aligned} H(X) &= \sum_{i=1}^n P(x_i) I(x_i) \\ &= - \sum_{i=1}^n P(x_i) \log P(x_i) \end{aligned}$$

- When the base of the logarithm is 2 the units of $I(x)$ are in bits
- The entropy of X can be interpreted as the expected value of

$$\log\left(\frac{1}{P(X)}\right)$$

- $H(X)$ is called the **Entropy**

Average Self Information

→ The term entropy has been borrowed from statistical mechanics. where it is used to denote the level of disorder in a system.

→ We observe that since $0 \leq P(x_i) \leq 1$

$$\log\left(\frac{1}{P(X)}\right) \geq 0$$

→ $H(X) \geq 0$

Example

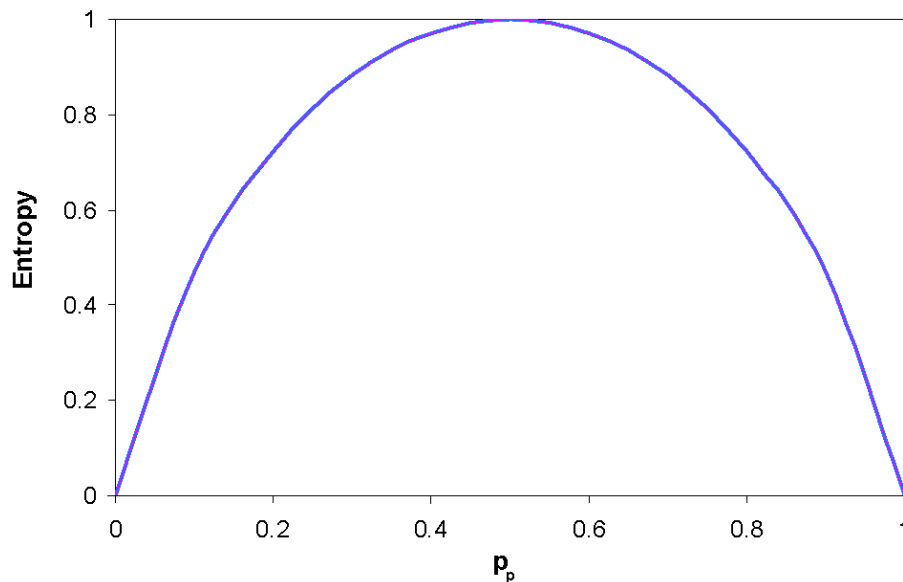
- ➔ Consider a discrete binary source that emits a sequence of statistically independent symbols.
- ➔ The output is either a 0 with probability p and a 1 with a probability $1 - p$.
- ➔ The entropy of this binary source is

$$H(X) = -(p)\log_2(p) - (1 - p)\log_2(1 - p)$$

Example

$$H(X) = -(p)\log_2(p) - (1 - p)\log_2(1 - p)$$

The plot of the binary entropy function vs p is



Entropy of English alphabet

- Consider the English language with alphabet {A. B. Z}.
- If every letter occurred with the same probability and was independent from the other letters, then the entropy per letter would be $\log_2 26 = 4.70$
- This is the absolute upper bound.
- However. we know that all letters do not appear with equal probability.
- S. T. A. E are more frequent
- Q. J. Z. J are less frequent

Entropy of English alphabet

- Consider the English language with alphabet $\{A. B. Z\}$.
- If we take into consideration the probabilities of occurrences of different alphabets (normalized letter frequency), the entropy per letter, H_L would be

$$H_L = 4.14bits \leq H(X)$$

- If X^2 denotes the random variable of bigrams in the English language. the upper bound on H_L can be defined as

$$H_L \leq H(X) \approx 3.56bits$$

- Here we consider the probabilities of all pairs.

Entropy of English alphabet

- The logic can be extended to n-grams. Thus the entropy of the language can be defined as

$$\lim_{n \rightarrow \infty} \frac{H(X^n)}{n}$$

- Even though the exact value of H_L is difficult to determine, statistical investigations show that for the English language $1 \leq H_L \leq 1.5 \text{ bits}$
- So each letter in the english text gives at most 1.5 bits of information.
- Let assume the value of H_L is 1.25 bits. Thus the redundancy of the English language is

$$\begin{aligned} R_{Eng} &= 1 - \frac{H_L}{\log_2 26} \\ &= 1 - \frac{1.25}{4.70} \approx 0.75 \end{aligned}$$

Entropy of Spoken English

- Let us now consider the redundancy in spoken english.
- Suppose an average speaker speaks **60 words per minute** and the average number of letters **per word is 6**.
- The average number of letters spoken per second in this case is **6 letters/sec**.
- Assuming **each letter carries 1.25 bits of information**, the information rate of an average speaker is **7.5 bits/sec**.
- If each letter is represented **by 5 bits**, the bitrate of an average speaker is **30 bit/sec**.
- **However, the typical data rate requirement for speech is 32 kilobits/sec.**

Conditional Entropy

→ The average conditional self information is called the conditional Entropy

$$H(X|Y) = \sum_{i=1}^n \sum_{j=1}^m P(x_i, y_j) \log\left(\frac{1}{P(x_i|y_j)}\right)$$

- The physical interpretation of this definition is as follows
- $H(X|Y)$ Is the information (or uncertainty) in X after Y is observed
- Based on the definition of $H(X|Y)$ we can write

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

Conditional Entropy and Mutual Information

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

- ➔ Since $I(X; Y) \geq 0$, it implies that $H(X) \geq H(X|Y)$
- ➔ The case $I(X; Y) = 0$ implies that $H(X) = H(X|Y)$ which is possible if and only if X and Y are statistically independent

Conditional Entropy and Mutual Information

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

- ➔ Since $H(X|Y)$ is the average amount of uncertainty(information) in X after we observe Y and $H(X)$ is the average amount of uncertainty (self information) of X. $I(X; Y)$ is the average amount of uncertainty (mutual information) about X having observed Y
- ➔ Since $H(X) \geq H(X|Y)$ the observation of Y does not increase the entropy(uncertainty). It can only decrease the entropy.

Joint Entropy

- The joint entropy of a pair of discrete random variables (X, Y) with a joint distribution $P(x, y)$ defined as

$$H(X, Y) = - \sum_{i=1}^n \sum_{j=1}^m P(x_i, y_j) \log P(x_i, y_j)$$

- By using the mathematical definitions of $H(X)$, $H(X, Y)$ and $H(X|Y)$ we obtain the following chain rule

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

- And consequently

$$I(X; Y) = H(X) + H(Y) - H(X, Y)$$

Venn Diagram

