

# Information Theory for Machine Learning

---

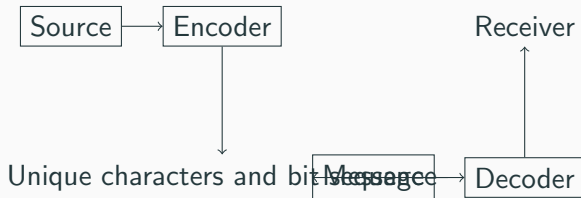
Nipun Batra

June 10, 2023

IIT Gandhinagar

# The Data Compression Problem

“hello”



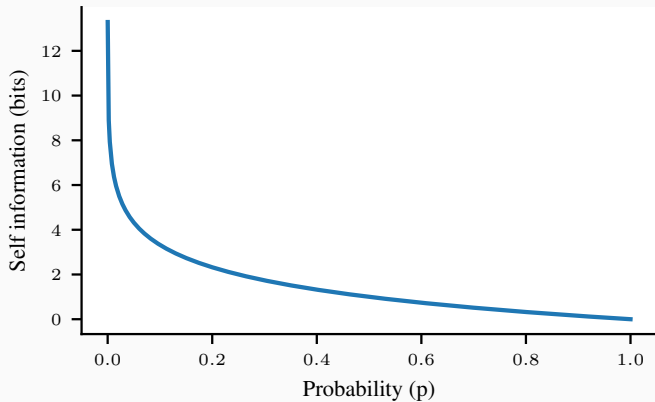
**Figure 1:** Data Compression Problem

- What is more surprising: Snowing in Kashmir or Snowing in Gandhinagar?
- To formalize, let us assume that the probability of snowing in Kashmir is  $p_1$  and that in Gandhinagar is  $p_2$ , and that  $p_1 \gg p_2$ .
- How can we quantify the surprise?

# Self Information

- Events that are less likely to occur are more surprising.
- Also, if an event is 100% likely to occur, it is not surprising at all.
- Also, if two events are independent, then the surprise of both of them occurring together is the sum of the surprise of each of them occurring individually.
- So, we need a function that maps probability to a number. Function should be: monotonic, and additive, and is 0 when the probability is 1.
- The function is  $I(x) = -\log_2(x)$  also called the self information or surprisal.

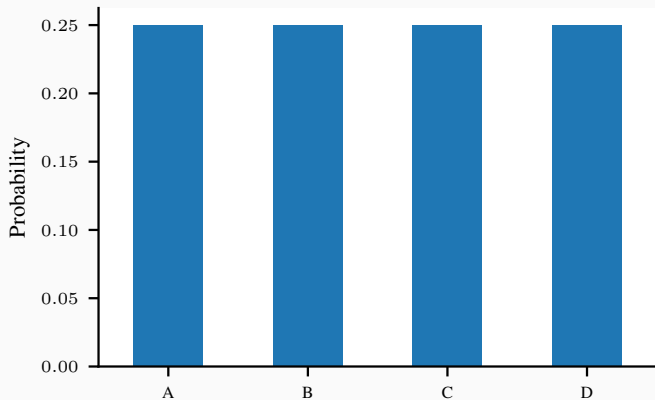
# Self Information



**Figure 2:** Self Information

## Self Information

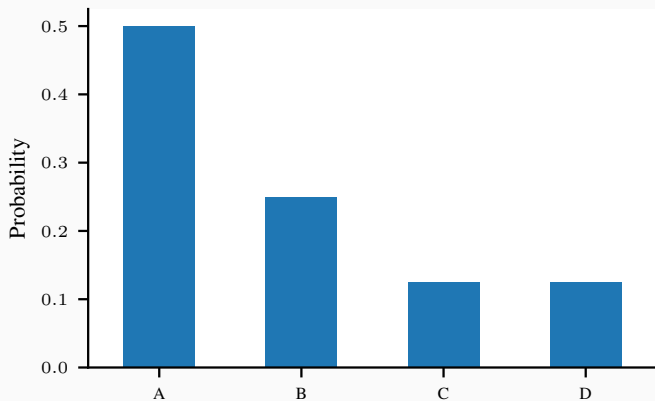
Consider a categorical random variable  $X$  with 4 possible outcomes: A, B, C, D. The probability of each of these outcomes is 0.25. What is the self information of each of these outcomes?



$$I(A) = I(B) = I(C) = I(D) = 2 \text{ bits.}$$

## Self Information

Consider a categorical random variable  $X$  with 4 possible outcomes: A, B, C, D. The probability these outcomes is 0.5, 0.25, 0.125, and 0.125. What is the self information of each of these outcomes?



$I(A) = 1$  bit,  $I(B) = 2$  bits,  $I(C) = I(D) = 3$  bits.

Proof on additivity of self information: Consider two independent random variables  $X$  and  $Y$  with PMFs  $p_X(x)$  and  $p_Y(y)$  respectively. The joint PMF is  $p_{X,Y}(x,y) = p_X(x)p_Y(y)$ . The self information of the joint PMF is:

$$\begin{aligned} I(X = x, Y = y) &= -\log_2(p_X(x)p_Y(y)) \\ &= -\log_2(p_X(x)) - \log_2(p_Y(y)) \\ &= I(X = x) + I(Y = y) \end{aligned}$$

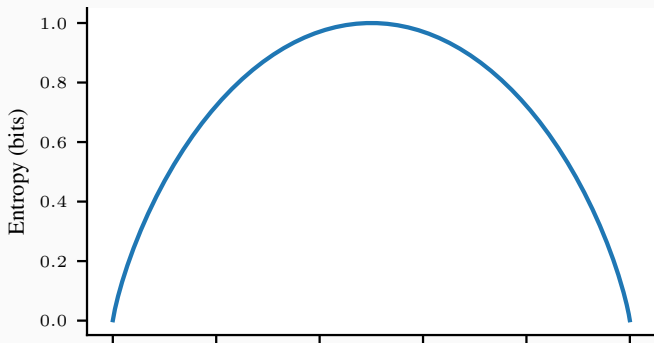


- The entropy of a random variable is the expected value of the self information.
- $H(X) = \mathbb{E}_{X \sim p(x)}[I(X)] = \mathbb{E}_{X \sim p(x)}[-\log_2(p(x))]$
- The entropy of a random variable is the expected number of bits required to encode the random variable.
- The entropy of a random variable is the minimum number of bits required to encode the random variable.

# Entropy

For a Bernoulli random variable  $X$  with probability  $p$  of success, the entropy is:

$$\begin{aligned} H(X) &= \mathbb{E}_{X \sim p(x)}[-\log_2(p(x))] \\ &= -\log_2(p) \times p - \log_2(1-p) \times (1-p) \\ &= -p \log_2(p) - (1-p) \log_2(1-p) \end{aligned}$$



# Entropy

For a  $k$  class categorical random variable  $X$  with probability  $p_i$  of class  $i$ , the entropy is:

$$\begin{aligned} H(X) &= \mathbb{E}_{X \sim p(x)}[-\log_2(p(x))] \\ &= -\sum_{i=1}^k p_i \log_2(p_i) \end{aligned}$$

