

Information theory

Kristian Wichmann

November 4, 2016

1 Self-information or surprisal

Let X be a random variable. Consider an event A . We may ask ourselves how much information $I(A)$ - also known as *self-information* or *surprisal* - we have gained by having this event occurring. It is clear, that such a quantity must depend only on the probability of the event:

$$I(A) = I(P(A)) \quad (1.1)$$

Therefore, we can express self-information through a function $f(p)$, so that if $P(A) = p$, then $I(A) = f(p)$.

If the outcome of an event A is certain, i.e. if $P(A) = 1$ then we have gained no information. So we must have $P(A) = 1 \Rightarrow I(A) = 0$. or in other words $f(1) = 0$. Non-certain events occurring, on the other hand, should give us non-zero information. So for $p < 1$ we should have $f(p) > 0$.

Further, if two events A and B are independent it seems reasonable to require that self-information is additive in the following sense:

$$I(A \cap B) = I(A) + I(B) \quad (1.2)$$

So if two independent events happen at the same time, self-information should simply add up. Because of independence, we also have:

$$P(A \cap B) = P(A) \cdot P(B) \quad (1.3)$$

Applying f to both sides of this equation we get:

$$I(A \cap B) = f(P(A) \cdot P(B)) \quad (1.4)$$

Combine this with equation (1.2) to get:

$$f(P(A) \cdot P(B)) = f(P(A)) + f(P(B)) \quad (1.5)$$

The only functions having this property are logarithms. Hence, the self-information must be of the form:

$$f(p) = -k \cdot \log(p) \quad (1.6)$$

The minus sign comes from requiring $f(p) > 0$ for $p < 1$. This means that k will be positive, but apart from that can be chosen freely. Since all logarithms are proportional to each other, this is equivalent to choice of base b being free:

$$f(p) = -\log_b(p) \quad (1.7)$$

1.1 Continuous distributions?

The section above deals with discrete random variables? However, we run into problems if we try to mindlessly generalize to continuous variables: The "obvious" analogue of the self-information for the outcome $X = x$ would be $-\log_b f(x)$, where $f(x)$ is the probability density function of X . But since this function need not be below 1, the associated surprisal may actually be negative! Clearly, something is fishy. But for now, we will only consider discrete random variables.

2 Entropy

The *entropy* of a discrete random variable X is the expectation value of the self-information:

$$H(X) = E[I(X)] = E[-\log_b(X)] \quad (2.1)$$

Here, $I(X)$ is itself a stochastic variable. Thus, entropy can be interpreted as the expected surprisal. Since X is discrete, we may write:

$$H(X) = - \sum_x p(x) \log_b p(x) \quad (2.2)$$

Figure 1 shows how much an outcome of p contributes to the total entropy. Since the limit for $p \rightarrow 0$ tends to zero, we will extend the definition to outcomes with zero probability; these do not contribute to the entropy.

2.1 Different choices of b

As mentioned above, we're free to choose b , but some choices are common. Each carry its own unit of entropy with it:

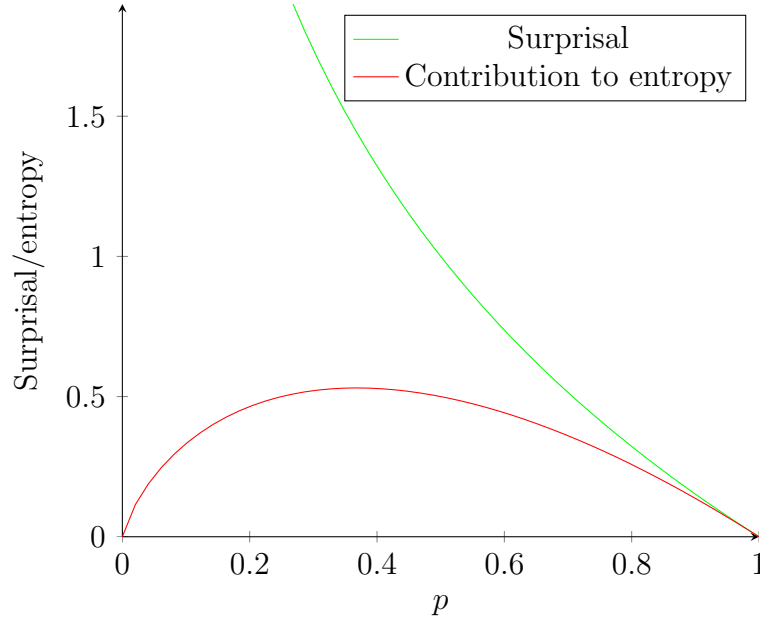


Figure 1: Surprisal and contribution to entropy as a function of p . Here for $b = 2$.

- $b = 2$: The corresponding entropy is known as *Shannon entropy*, and the unit is Shannon or simply bits.
- $b = e$: The corresponding unit is known as a nat.
- $b = 10$: The corresponding unit is known as a Hartley.

Unless explicitly mentioned, we will use Shannon entropy from now on.

2.2 Example: Entropy of a coin toss

Let's consider the simplest possible non-trivial situation: an experiment with two outcomes, A and B . If the probability of A is p , then the probability of B must be $1 - p$. For a fair coin, both probabilities are $\frac{1}{2}$, and the entropy is:

$$H = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = \frac{1}{2} + \frac{1}{2} = 1 \quad (2.3)$$

If the coin is not fair, but the probability of tails (A) is p , instead we get:

$$H(p) = -p \log_2 p - (1 - p) \log_2 (1 - p) \quad (2.4)$$

This function is plotted in figure 2.

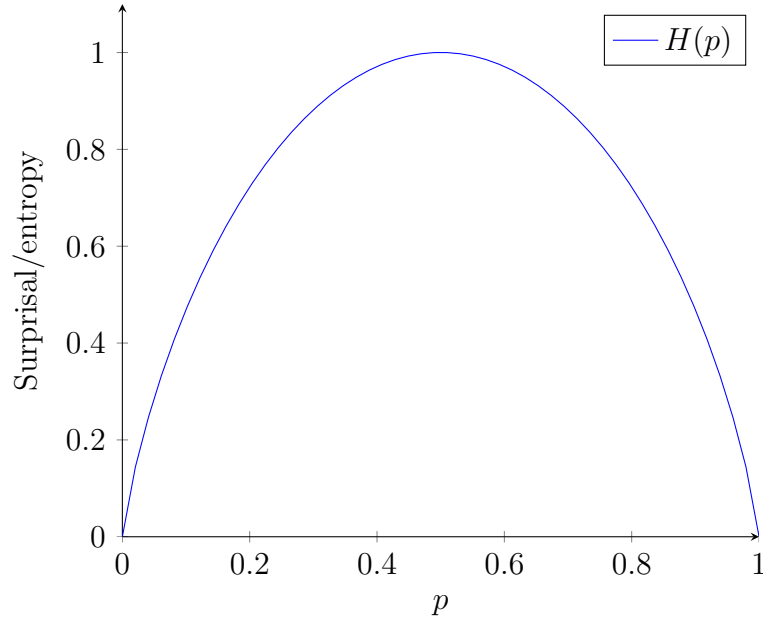


Figure 2: Entropy for an unfair coin toss.

3 Conditional entropy and mutual information

3.1 Entropy of a joint distribution

Given two discrete random variables X and Y , we may define their *joint entropy* simply as the entropy of their joint distribution:

$$H(X, Y) = - \sum_{x,y} P(X = x, Y = y) \log_2 (P(X = x, Y = y)) \quad (3.1)$$

Or using the joint distribution function $p(x, y) = P(X = x, Y = y)$:

$$H(X, Y) = - \sum_{x,y} p(x, y) \log_2 p(x, y) \quad (3.2)$$

3.2 Conditional entropy

Similarly, we may define *conditional entropy*. If we already know the outcome of one random variable, this will limit the number of outcomes that contributes to the entropy. But the probabilities become conditional:

$$H(X|Y = y) = - \sum_x p(x|y) \log_2 p(x|y) \quad (3.3)$$

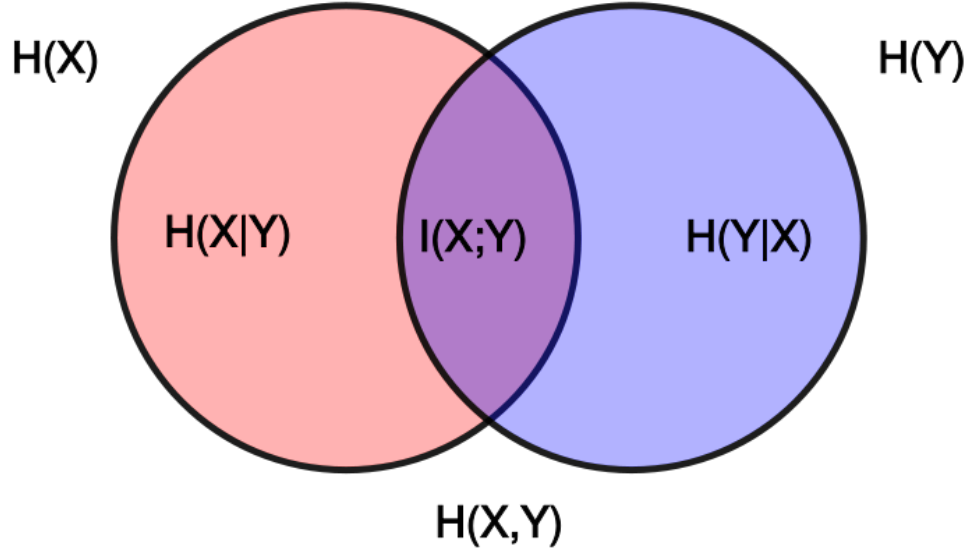


Figure 3: Visualization of the different entropies and mutual information. Image source: Wikipedia.

Since $p(x|y) = \frac{p(x,y)}{p(y)}$ this means:

$$H(X|Y = y) = - \sum_x \frac{p(x,y)}{p(y)} \log_2 \frac{p(x,y)}{p(y)} \quad (3.4)$$

The total conditional entropy is found by weighing all of these:

$$H(X|Y) = \sum_y p(y) \cdot H(X|Y = y) \quad (3.5)$$

Here $p(y) = \sum_x p(x,y)$ is the marginal probability for Y . (Similarly $p(x) = \sum_y p(x,y)$). Inserting equation 3.4 into equation 3.5 we get:

$$H(X|Y) = - \sum_y p(y) \sum_x \frac{p(x,y)}{p(y)} \log_2 \frac{p(x,y)}{p(y)} = - \sum_{xy} p(x,y) \log_2 \frac{p(x,y)}{p(y)} \quad (3.6)$$

Looking at figure 3, we would expect that $H(X,Y) - H(Y) = H(X|Y)$.

Let's check that this is indeed the case:

$$H(X, Y) - H(Y) = - \sum_{x,y} p(x, y) \log_2 p(x, y) - \left(- \sum_y p(y) \log_2 p(y) \right) \quad (3.7)$$

Use the definition of $p(y)$:

$$\sum_y p(y) \log_2 p(y) = \sum_{xy} p(x, y) \log_2 p(y) \quad (3.8)$$

Hence:

$$H(X, Y) - H(Y) = - \sum_{x,y} p(x, y) (\log_2 p(x, y) - \log_2 p(y)) \quad (3.9)$$

This is the same as:

$$- \sum_{x,y} p(x, y) \log_2 \frac{p(x, y)}{p(y)} = H(X|Y) \quad (3.10)$$

3.3 Mutual information

If we look at figure 3 once again, we also notice $I(X; Y)$, the *mutual information* between two discrete random variables X and Y . It should be equal to:

$$I(X; Y) = H(X) - H(X|Y) \quad (3.11)$$

Inserting we get:

$$- \sum_x p(x) \log_2 p(x) - \left(- \sum_{x,y} p(x, y) \log_2 \frac{p(x, y)}{p(y)} \right) \quad (3.12)$$

Now use $p(x) = \sum_y p(x, y)$ to rewrite the first term:

$$\sum_x p(x) \log_2 p(x) = \sum_{xy} p(x, y) \log_2 p(x) \quad (3.13)$$

Combine terms to get:

$$I(X; Y) = \sum_{x,y} p(x, y) \log_2 \frac{p(x, y)}{p(x)p(y)} \quad (3.14)$$

This is clearly symmetrical: $I(X; Y) = I(Y; X)$. Also, if X and Y are independent, then $p(x, y) = p(x)p(y)$ which makes the fraction 1 and logarithm zero, so the mutual information vanishes in this case.