Exercise: We derive the Shannon entropy as a measure of the uncertainty represented by a discrete probability distribution.

(Source: Jaynes, Information Theory and Statistical Mechanics)



Suppose we have a variable $X$ which can take values $\{X_1, \ldots, X_n\}$. Our incomplete understanding of the processes which determine the value of $X$ can be represented by assigning probability $p_i$ to outcome $X_i$, for $i \in \{1, \ldots, N\}$

We wish to find a function $H$ to characterize the "uncertainty" represented by the probability distribution $\{p_i\}_{i=1}^{n}$. We require $H$ to have 3 properties

1.) $H(p_1, \ldots, p_n)$ is continuous

2.) When all outcomes are equally probable (i.e. $p_i = 1/n \ \forall i \in \{1, \ldots, n\}$),

$$H(1/n, \ldots, 1/n) \equiv A(n)$$ is a monotonically increasing function of $n$.

(Intuitively, a uniform distribution over more outcomes is more "uncertain" than

a uniform distribution over fewer values.)

## 3.) Composition Law

Suppose we wish to group outcomes together. We might group $\{X_1, \ldots, X_k\}$ together with probability $W_1 = P_1 + \cdots + P_k$, $\{X_{k+1}, \ldots, X_{k+m}\}$ with probability $W_2 = P_{k+1} + \cdots + P_{k+m}$, and so on. We can then assign conditional probabilities as $P(X_1 | W_1) = P_1 / W_1, \ldots, P(X_k | W_1) = P_k / W_1$, and so on. Specifying the group probabilities and the conditional probabilities is equivalent to specifying the original probabilities, so we require the uncertainty to be the same in both cases.

$$H(P_1, \ldots, P_n) = H(W_1, \ldots, W_r) + W_1 H\left(\frac{P_1}{W_1}, \ldots, \frac{P_k}{W_1}\right)$$

$$+ W_2 H\left(\frac{P_{k+1}}{W_2}, \ldots, \frac{P_{k+m}}{W_2}\right) + \cdots$$

We show that the Shannon entropy satisfies these requirements.

## Proof

Suppose we have probabilities $\{P_i\}_{i=1}^n$. By (1), we only need to consider rational probabilities, as irrational numbers can be constructed from

sequences of rational numbers, and this continuous, so it will behave well under such sequence limits. Thus, we can find numbers $\{n_i\}_{i=1}^N$, $n_i \in \mathbb{N}$ such that $P_i = n_i / \sum\limits_i n_i$ .

We will now treat these probabilities as "groups" made from some uniform distribution over $\sum\limits_i n_i$ numbers. By rule (2), this uniform distribution has uncertainty $A(\sum\limits_i n_i)$.

By rule (3), we have

$$A\left(\sum_i n_i\right) = H(P_1, \ldots, P_n) +$$

$$\sum_j P_i \, H\left(\underbrace{\frac{1/\sum\limits_i n_i}{n_j / \sum\limits_i n_i}, \ldots, \frac{1/\sum\limits_i n_i}{n_j / \sum\limits_i n_i}}_{A(n_j)}\right) \Longrightarrow$$

$$H(P_1, \ldots, P_n) = A\left(\sum_i n_i\right) - \sum_i P_i \, A(n_i).$$

To determine $A$, we consider the case when $n_i = m$ $\forall i \in \{1, \ldots, n\}$.

Then, $\sum_i n_i = m \cdot n$, and

$$H\left(\frac{m}{mn}, \cdots, \frac{m}{mn}\right) = A(mn) - \sum_i \frac{m}{mn} A(m)$$

$$\Longrightarrow \quad A(n) + A(m) = A(m \cdot n)$$

We thus conclude that $A(n) = K \log(n)$, where $K > 0$ by rule (2). Finally, substituting into the formula for $H(p_1, \cdots, p_n)$, we derive

$$H(p_1, \cdots, p_n) = K \log\left(\sum_i n_i\right) - K \sum_i p_i \log(n_i) =$$

$$K \log\left(\sum_i n_i\right) - K \sum_i \left[p_i \log(n_i) - p_i \log\left(\sum_i n_i\right) + p_i \log\left(\sum_i n_i\right)\right]$$

$$= K \log\left(\sum_i n_i\right) - K \sum_i p_i \log\left(n_i / \sum_i n_i\right) - K \log\left(\sum_i n_i\right) \sum_i p_i$$

$$= -K \sum_i p_i \log p_i \qquad \Longleftrightarrow$$

$$\boxed{H(p_1, \cdots, p_n) = -K \sum_{i=1}^{n} p_i \log p_i}$$