

Day 13 - Properties of the Entropy

Elements of Information Theory

Definition 1. The **entropy** of an event A is a measure of the uncertainty we feel about the occurrence of A .

The entropy of a random variable X is given by

$$H(X) = \sum_a \mathbb{P}(X=a) \cdot \log_2 \left(\frac{1}{\mathbb{P}(X=a)} \right)$$

all outcomes

amount of info we gain by knowing $X=a$.

Entropy = uncertainty of a random variable.
= average amount of information we gain by knowing the outcome of a R.V.

= average number of bits we need to store the outcomes of R.V

Definition 2. The entropy of two random variables X and Y is

$$H(X,Y) = \sum_{a,b} \mathbb{P}(X=a, Y=b) \cdot \log_2 \left(\frac{1}{\mathbb{P}(X=a, Y=b)} \right)$$

all outcomes $(X,Y)=(a,b)$ $\mathbb{P}(X=a \cap Y=b)$

the average amount of info we gain by learning the outcome of both X & Y

Definition 3. The conditional entropy of the random variable X given an event B is

$$H(X|B) = \sum_a \mathbb{P}(X=a|B) \cdot \log_2 \left(\frac{1}{\mathbb{P}(X=a|B)} \right)$$

The additional amount of info we gain by learning X after we know that B has occurred.

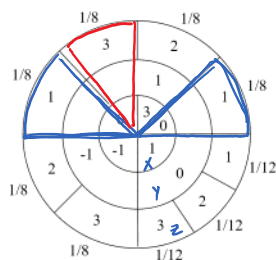
Definition 4. The conditional entropy of the random variable X given a random variable Y is

$$H(X|Y) = \sum_b \mathbb{P}(Y=b) \cdot H(X|Y=b)$$

sum over all the condition events

$\mathbb{P}(\text{condition event happens})$

Example. Suppose that random variables X, Y, Z are obtained by spinning the wheel below, with X given by the innermost circle, Y given by the intermediate circle, and Z given by the outermost circle.



last time

- (a) Compute $H(X)$
- (b) How many bits (of information) are required to store the results of 100,000 spins of Z ? $= 100,000 \cdot H(2)$
- (c) Calculate the uncertainty of Z given that $X = 0$.

$$\begin{aligned}
 H(Z|X=0) &= \sum_a P(Z=a | X=0) \cdot \log_2 \left(\frac{1}{P(Z=a | X=0)} \right) \\
 &= \left(\frac{2}{3} \right) \cdot \log_2 \left(\frac{1}{2/3} \right) + \left(\frac{1}{3} \right) \cdot \log_2 \left(\frac{1}{1/3} \right) \\
 &\approx 0.9183
 \end{aligned}$$

$$P(Z=1 \cap X=0) = \frac{2}{8}$$

$$P(Z=2 \cap X=0) = 0$$

↳ not on the wheel.

$$P(Z=3 \cap X=0) = \frac{1}{8}$$

$$P(Z=1 | X=0) = \frac{P(Z=1 \cap X=0)}{P(X=0)} = \frac{2/8}{3/8} = \frac{2}{3}$$

$$P(Z=2 | X=0) = 0$$

$$P(Z=3 | X=0) = \frac{P(Z=3 \cap X=0)}{P(X=0)} = \frac{1/8}{3/8} = \frac{1}{3}$$

(d) Calculate $H(Z|Y) = \sum_y P(Y=y) \cdot H(Z|Y=y)$

One way to do this: just compute all $H(Z|Y=-1)$, $H(Z|Y=0)$ and $H(Z|Y=1)$.

A faster way is to realize that Z and Y are independent random variables

Defn: 2 R.V. X and Y are independent if & only if.

$$P(X=a \cap Y=b) = P(X=a) \cdot P(Y=b) (*)$$

for all values $(x, y) = (a, b)$.

You check that, under this wheel, Z & Y are independent - by checking $(*)$ holds for all pairs of values for (Z, Y)

Let X, Y be random variables. Then

$$H(Y|X) = H(Y) \Leftrightarrow X \text{ and } Y \text{ are independent.}$$

According to this, $H(Z|Y) = H(Z) = \log_2(3)$

(e) Calculate $H(X|Y, Z)$
 Let X, Y_1, Y_2, \dots, Y_k be random variables then
 $H(X|Y_1, Y_2, \dots, Y_k) = 0 \Leftrightarrow X = f(Y_1, \dots, Y_k)$
 You don't gain any additional info for knowing X after knowing all Y_i 's.
 X is a fcn of Y_1, Y_2, \dots, Y_k .

$H(X|Y, Z)$

Try to see if you can write X as a fcn of Y and Z .

If we have, for example: $X = Y^2 + 5Z^3$
 then $H(X|Y, Z) = 0$.

Suppose we learn the value of X and then we learn the value of Y . Then

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

info about X, Y = info about X + remaining info that we still don't know about Y , after knowing X

$$H(x_1, x_2, \dots, x_N) = H(x_1) + H(x_2|x_1) + H(x_3|x_1, x_2) + \dots + H(x_N|x_1, \dots, x_{N-1})$$

An observation: $P(X=a \cap Y=b) = P(X=a) \cdot \frac{P(X=a \cap Y=b)}{P(X=a)} = P(X=a) \cdot P(Y=b|X=a)$

Proof: $H(X, Y) = \sum_a \sum_b P(X=a \cap Y=b) \cdot \log_2 \left(\frac{1}{P(X=a \cap Y=b)} \right)$

$$= \sum_a \sum_b P(X=a \cap Y=b) \cdot \log_2 \left(\frac{1}{P(X=a) \cdot P(Y=b|X=a)} \right)$$

$$= \sum_a \sum_b P(X=a \cap Y=b) \cdot \log_2 \left(\frac{1}{P(X=a)} \right) + \sum_a \sum_b P(X=a \cap Y=b) \cdot \log_2 \left(\frac{1}{P(Y=b|X=a)} \right)$$

$$= \sum_a \sum_b P(X=a) \cdot \log_2 \left(\frac{1}{P(X=a)} \right) \cdot P(Y=b|X=a) + \sum_a \sum_b P(X=a) \cdot P(Y=b|X=a) \cdot \log_2 \left(\frac{1}{P(Y=b|X=a)} \right)$$

$$= \underbrace{\sum_a P(X=a) \cdot \log_2 \left(\frac{1}{P(X=a)} \right)}_{H(X)} \cdot \underbrace{\sum_b P(Y=b|X=a)}_1 + \sum_a P(X=a) \cdot \underbrace{H(Y|X=a)}_{H(Y|X)}$$

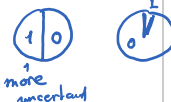
Important inequalities

1. For a random variable X which takes only k values we always have

$$H(X) \leq \log_2(k)$$

with equality if and only if X takes all its values with equal probability

max uncertainty about X occurs when all its values are equally likely.



2. For any two random variables X and Y we always have

$$H(X|Y) \leq H(X)$$

and equality holds if and only if X and Y are independent

the amount of info we gain by learning X after knowing Y
is less than the amount of info
we would gain if we did not
know Y .

3. For any two random variables X and Y we always have

$$H(X, Y) \leq H(X) + H(Y)$$

and equality holds if and only if X and Y are independent.

the amount of info. we gain by learning X & Y is less
than the amount of info we gain if.
we learn these X, Y separately