# INFormation Theory Foundations and Applications

# NOTe 1

# INTRODUtion to information theory

# Shannon Theory

## (A) INTRODUCTION

Information theory answers two fundamental questions:

- What is the ultimate data compression?

    Answer: Shannon entropy!

- What is the ultimate transmission rate of communication?

    Answer: The channel capacity!

These problems were solved by what is known today as Shannon coding theorems, that are the basis of information theory.

- Statistical physics (thermodynamics)
- Computer science (Kolmogorov complexity)
- Statistical inference (Occam's razor)
- Probability theory (hypothesis testing)

Are some developments that follows from information theory.

Before proving the theorems, we need to introduce some fundamental concepts, starting with the very notion of information.

① What is information and how to measure it

One of the fundamental contributions of Shannon is the notion of a bit as a measure of information.

Physical bit
⇓
two states of a physical system
(magnetic spin)

Shannon bit
⇓
measure of the surprise upon learning the outcome of a random binary experiment

If we toss a fair coin and look at the outcome, we learn one bit of information.

The outcome of a coin flip is the physical bit, but it is the information associated with the random nature of the physical bit that we want to measure.

Now that we have a unit to measure information, we need to define the measure.

✳ Let us assume that every physical system can be described as a random variable

$$X = \{ P_X(x), x \in \mathcal{X} \}$$

$\mathcal{X} \rightarrow$ alphabet

$x \rightarrow$ realization of the random variable

$P_X(x) \rightarrow$ probability distribution associated to $x$

Shannon's notion of information contained in the occurrence of an event.

i) I must be a function only on the probability

ii) I must be a continuous function

iii) I must be additive for independent events.

There is only one function that respect there postulates.

Let us consider $s$ independent occurrences of the event $x$. Then, $I$ must be a function of the total probability $[P_x(x)]^s$.

$$I\left([P_x(x)]^s\right) \stackrel{\text{Ind.}}{=\!=} I\left([P_x(x)]^{s-1}, P_x(x)\right)$$

$$\stackrel{\text{Add.}}{=\!=} I\left([P_x(x)]^{s-1}\right) + I\left(P_x(x)\right)$$

$$\stackrel{\text{2nd.}}{=\!=} I\left([P_x(x)]^{s-2}, P_x(x)\right) + I\left(P_x(x)\right)$$

$$\stackrel{\text{Add.}}{=\!=} I\left([P_x(x)]^{s-2}\right) + 2\, I\left(P_x(x)\right)$$

$$\vdots$$

$$= s\, I\left(P_x(x)\right)$$

As a consequence, for any integer $t$ we have

$$I\left([P_x(x)]^{1/t}\right) = \frac{t}{t}\, I\left([P_x(x)]^{1/t}\right) = \frac{1}{t}\, I\left([P_x(x)]^{t/t}\right)$$

$$= \frac{t}{t}\, I\left(P_x(x)\right)$$

Therefore, for any rational number

$$r = \frac{s}{t}$$

we must have

$$I\left([p_x(x)]^r\right) = r\, I\left(P_x(x)\right)$$

Now, any probability can be written as

$$P_x(x) = 2^{\log P_x(x)}$$

And any real number can be arbitrarily well approximated by a rational number. Then

$$I\left(P_x(x)\right) = I\left(2^{\log P_x(x)}\right) = \log P_x(x)\, I(2)$$

So, we choose $I(2) = -1$ to get

$$I(p_x(x)) = -\log P_x(x)$$

This is the amount of information contained in the event $x$. It is how much we learn from knowing the value of $X$.

I is the measure of the information contained in a single occurrence of the random variable. We are interested in the information contained in the physical system, which is the information source. Therefore, we define the average information

$$H(X) = - \sum_{x \in \chi} P_X(x) \log P_X(x) = \mathbb{E}\left[ I(P_X(x)) \right]$$

That is Shannon entropy, which measures the uncertaint we have about $X$, or how much information we gain when we learn the value of $X$.

For a fair coin we have

$$P_X(x) = \left\{ \frac{1}{2}, \frac{1}{2} \right\}$$

$$x = \{ H, T \}$$

$$\Rightarrow$$

$$H(X) = - \frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2}$$

$$= 1 \text{ bit}$$

② Some properties of Shannon entropy

a) Entropy is non-negative

$$H(x) \geq 0$$

This follows because it is the average of a positive quantity.

b) The entropy is invariant with respect to the permutations of the realizations of $X$.

This is because it depends only on the probabilities, not on the values of the realizations

c) $H(x) = 0$ for a deterministic variable

Let us consider a deterministic distribution

$$P_x(x) = \delta_{x,x_0}$$

$$\Rightarrow H(x) = 0$$

$\Leftarrow$ If $H(x) = 0$, then we have

$P_X(x) \log \frac{1}{P_X(x)} = 0$ for all $x \in X$, which

implies $P_X(x) = 0$ or $P_X(x) = 1$. Since

$P_X(x)$ must be a probability distribution,

we must have $P_X(x_0) = 1$ and $P_X(x) = 0$ for

all others values of $x$.

This is intuitively expected from the

meaning of entropy.

d) $H(x)$ is upper bounded

$$H(x) \leq \log |x|$$

with $|x|$ being the cardinality of $X$.

First, let us consider a uniform random variable

$$P_X(x) = \frac{1}{|X|} \quad \forall x$$

For this case we have

$$H(x) = \log |x|$$

Let us now move to the inequality.

We consider a Lagrangian optimization[1], with the Lagrangian being defined as

$$\mathcal{L} = H(x) + \lambda \left( \sum_x P_x(x) - 1 \right)$$

$$\delta \mathcal{L} = \left\{ - \sum_x \left[ \log(P_x(x)) - 1 \right] + \lambda \sum_x 1 \right\} \delta P_x(x) = 0$$

$$\Rightarrow \quad -\log P_x - 1 + \lambda = 0 \quad \Rightarrow \quad P_x(x) = 2^{\lambda - 1}$$

Since $\lambda$ is constant, the probability distribution that maximizes $H(x)$ is the uniform one.

Therefore we conclude that

$$\boxed{0 \leq H(x) \leq \log |x|}$$

[1] We assume that the entropy is concave. We will prove this latter.

③ Other measures of information

Ⓐ Conditional entropy

If two random variables are correlated, by measuring one of them we obtain information about the other.

Let us define the conditional information content

$$i(x|y) = -\log\left(P_{X|Y}(x|y)\right)$$

The conditional entropy is defined as the expected conditional information content

$$H(X|Y) = \mathbb{E}_{X,Y}\left\{i(X|Y)\right\} = \sum_{y}' P_Y(y)\, H(X|Y=y)$$

$$= -\sum_{x,y}' P_{X,Y}(x,y)\log P_{X|Y}(x|y)$$

Where we used $P_{X,Y} = P_Y(y)\, P_{X|Y}(x,y)$

$H(X|Y)$ is the amount of uncertainty about $X$ when we know $Y$.

(B)  Joint Entropy

It is the entropy of the joint random variable $(X, Y)$

$$H(X,Y) = \mathbb{E}_{X,Y}\{i(X,Y)\} = -\sum_{x,y} P_{X,Y}(x,y) \log P_{X,Y}(x,y)$$

$$H(X,Y) = -\sum_{x,y} P_{X,Y}(x,y) \left\{ \log P_X(x) + \log P_{Y|X}(y|x) \right\}$$

$$= -\sum_{x,y} P_{X,Y}(x,y) \log P_X(x) - \sum_{x,y} P_{X,Y}(x,y) \log P_{Y|X}(y|x)$$

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

Entropy is subadditive

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

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

$$\Rightarrow \boxed{H(X,Y) \leq H(X) + H(Y)}$$

© The mutual information

It is a measure of the correlations between two random variables

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

$$= \sum_{x,y} P_{X,Y}(x,y) \log\left(\frac{P_{X,Y}(x,y)}{P_X(x)\, P_Y(y)}\right)$$

Since $H(X) \geq H(X|Y) \Rightarrow \boxed{I(X:Y) \geq 0}$

The equality is achieved if and only if the random variables are independent

$$P_{X,Y}(x,y) = P_X(x)\, P_Y(y) \Rightarrow I(X:Y) = 0$$

# D) Relative entropy

The relative entropy is a measure of how far one probability distribution $P_X(x)$ is from another one $q_X(x)$.

It is defined as

$$D(P \| q) = \sum_x P_X(x) \log\left(\frac{P_X(x)}{q_X(x)}\right)$$

If $\text{supp}(P) \subseteq \text{supp}(q)$

The mutual information can be written in terms of the relative entropy as

$$I(X:Y) = D\left(P_{X,Y}(x,y) \| P_X(x) \times P_Y(y)\right)$$

This tells us how far we are from Independence!

We can use this result to prove that the entropy IS A concave function.

$$H(X) = -\sum_x' P_x(x) \log P_x(x) = -\sum_x P_x \log\left(\frac{P_x U_x}{U_x}\right)$$

$$= -\sum_x P_x \log \frac{P_x}{U_x} - \sum_x P_x \log U_x$$

Taking $U_x = \frac{1}{|X|}$ (uniform distribution)

$$H(X) = -D(P_x \| U_x) + \log |X|$$

$$\Rightarrow \log|X| - H(X) = D(P_x \| U_x)$$

Now, we have that $D(P_x \| q_x)$ is convex

$$D(\lambda P_1 + (1-\lambda)P_2 \| \lambda q_1 + (1-\lambda) q_2) = \sum_x' [\lambda P_1 + (1-\lambda) P_2] \log \frac{\lambda P_1 + (1-\lambda) P_2}{\lambda q_1 + (1-\lambda) q_2}$$

$$\leq \sum_x \left[ \lambda P_1 \log \frac{\lambda P_1}{\lambda q_1} + (1-\lambda) P_2 \log \frac{(1-\lambda) P_2}{(1-\lambda) q_2} \right]$$

$$= \lambda D(P_1 \| q_1) + (1-\lambda) D(P_2 \| q_2) \qquad \square$$

The second step Follows from the fact that, for real positive numbers $a_i$ and $b_i$

$$\left(\sum_{i=1}^{n} a_n\right) \log \frac{\sum_i a_v}{\sum_i b_i} \leq \sum_i a_v \log \frac{a_v}{b_i}$$

Therefore, for our special case

$$D\left(\lambda P_a + (1-\lambda)P_2 \| \lambda U_1 + (1-\lambda)U_2\right) \leq \lambda D(P_a \| U_1) + (1-\lambda)D(P_2 \| U_2)$$

$$\Rightarrow D\left(\lambda P_1 + (1-\lambda)P_2 \| U\right) \leq \lambda D(P_1 \| U) + (1-\lambda)D(P_2 \| U)$$

which implies

$$\log|X| - H\left(\lambda P_1 + (1-\lambda)P_2\right) \leq \lambda\left(\log|X| - H(P_1)\right)$$
$$+ (1-\lambda)\left(\log|X| - H(P_2)\right)$$
$$= \log|X| - \lambda H(P_1) - (1-\lambda)H(P_2)$$

$$\Rightarrow H\left(\lambda P_1 + (1-\lambda)P_2\right) \geq \lambda H(P_1) + (1-\lambda)H(P_2)$$

$$\boxed{H(X) \text{ is concave}}$$

# 4 Data processing inequality

Let $p$ and $q$ be two probability distributions, and let $\Delta$ be a classical channel. Then

$$D(p \| q) \geq D(\Delta p \| \Delta q)$$

## Proof

If $\text{supp}(p) \not\subseteq \text{supp}(q)$, then $D(p \| q) = \infty$ and the inequality is trivial.

If $\text{supp}(p) \subseteq \text{supp}(q)$, we have

$$\text{supp}(\Delta p) \subseteq \text{supp}(\Delta q)$$

Let us start by rewriting the the quantities appearing in the inequality.

$$D(\Delta p \| \Delta q) = \sum_y (\Delta p)(y) \log \frac{(\Delta p)(y)}{(\Delta q)(y)}$$

$$D(\Delta p \| \Delta q) = \sum_{x,y} \Delta(y|x) p(x) \log \frac{(\Delta p)(y)}{(\Delta q)(x)}$$

$$= \sum_{x} p(x) \left[ \sum_{y} \Delta(y|x) \log \frac{\Delta p)(y)}{(\Delta q)(y)} \right]$$

$$= \sum_{x} p_x \log \exp \left[ \sum_{y} \Delta(y|x) \log \frac{\Delta p(y)}{\Delta q(y)} \right]$$

$\Delta(y|x)$ is the conditional probability distribution defining the classical channel

$$X \xrightarrow{\Delta} Y$$

This implies that

$$D(p\|q) - D(\Delta p \| \Delta q) = D(p\|r)$$

$$r = q(x) \exp \left[ \sum_{y} \Delta(y|x) \log \frac{\Delta p(y)}{\Delta q(y)} \right]$$

Now, note that

$$\sum_x r(x) = \sum_x q \exp\left\{ \sum_y \Delta(y|x) \log \frac{\Delta p(y)}{\Delta q(y)} \right\}$$

$$\leq \sum_x q(x) \sum_y \Delta(y|x) \exp\left\{ \log \frac{\Delta p(y)}{\Delta \tilde{q}(y)} \right\}$$

$$= \sum_x q(x) \sum_y \Delta(y|x) \frac{\Delta p(y)}{\Delta \tilde{q}(y)}$$

$$= \sum_y \left[ \sum_x \tilde{q}(x) \Delta(y|x) \right] \frac{\Delta p(y)}{\Delta \tilde{q}(y)}$$

$$= \sum_y \Delta p(y) = 1$$

$$\Rightarrow \sum_x r(x) \leq 1$$

$$\Rightarrow D(p\|r) \geq 0$$

which proves Data Processing Inequality!

⑤ Fano's Inequality

$$X \xrightarrow{P_{Y|X}(y|x)} Y$$

↳ noisy communication channel

Y is processed and the best estimator $\hat{X}$ of X is produced. The probability error is

$$P_e = Pr\{\hat{X} \neq X\}$$

If the channel is noiseless, we have

$$P_{Y|X}(y|x) = \delta_{y,x} \Rightarrow H(X|Y) = 0$$

If the noise increases, $H(X|Y)$ increases

$H(X|Y)$ quantifies the amount of information lost in the channel.

Fano's inequality provides a quantitative relation between Pe and $H(X|Y)$

let us assume

$$X \longrightarrow Y \longrightarrow \hat{X}$$

Then

$$H(X|Y) \leq H(X|\hat{X}) \leq h_2(P_e) + P_e \log\left(|X|-1\right)$$

with $h_2(p) = -p\log p - (1-p)\log(1-p)$

Note that

$$\lim_{P_e \to 0} \left(h_2(P_e) + P_e \log\left(|X|-1\right)\right) = 0$$

$$\implies H(X|Y) = 0$$

As It should.

# Proof:

Let $E$ denote an error indicator

$$E = \begin{cases} 0 : & X = \hat{X} \\ 1 : & X \neq \hat{X} \end{cases}$$

Consider the entropy

$$H\left(EX \mid \hat{X}\right) = H\left(X \mid \hat{X}\right) + H\left(E \mid X\hat{X}\right)$$

If we know both $X$ and $\hat{X}$, there is no uncertainty about $E$. Therefore

$$H\left(E \mid X\hat{X}\right) = 0$$

And

$$H\left(EX \mid \hat{X}\right) = H\left(X \mid \hat{X}\right) \qquad \boxed{1}$$

Now, let us consider the following chain

$$X \rightarrow Y \rightarrow \hat{X}$$

then, we have

$$I(X:Y) \geq I(X:\hat{X}) \Rightarrow H(X|\hat{X}) \geq H(X|Y) \quad \textcircled{2}$$

$$\hookrightarrow \text{DATA processing inequality}$$

Now, we have

$$H(EX|\hat{X}) = H(E|\hat{X}) + H(X|E\hat{X})$$

Conditioning
Reduces entropy $\leq H(E) + H(X|E\hat{X})$

$$= h_2(p_e) + p_e H(X|\hat{X}, E=1)$$

$$+ (1-p_e) H(X|\hat{X}, E=0)$$

$$\leq h_2(p_e) + p_e \log(|X|-1) \quad \textcircled{3}$$

When there is no error $(E=0)$, there is
no uncertainty about $X$. Also, the uncertainty
about $X$, when $\hat{X}$ is available and we have an
error $(E=1)$, is less than the uncertainty
of a uniform distribution $1/(|X|-1)$. FANO'S
inequality follows from $\textcircled{1}$, $\textcircled{2}$ and $\textcircled{3}$.