

Elements of Information Theory

Sheng Yang
sheng.yang@centralesupelec.fr

“The fundamental problem of communication is that of reproducing at one point either exactly or approximately a message selected at another point.”

“A mathematical theory of communication”, 1948, Claude Shannon (1916-2001)

References

- T. Cover and J. Thomas, “Elements of information theory”
- Y. Polyanskiy and Y. Wu, “Information theory”
- I. Csiszár and J. Körner, “Information theory: Coding theorems for discrete memoryless systems”
- R. Gallager, “Information theory and reliable communication”
- R. Yeung, “A first course in information theory”
- A. El Gamal and Y.-H. Kim, “Network information theory”

Notations and terms

Throughout this course, we use the following notations and terminologies.

$\mathbb{R}, \mathbb{C}, \mathbb{Z}, \mathbb{N}$	real, complex, integer, natural numbers
i	$\sqrt{-1}$
x^n	(x_1, \dots, x_n)
$ \mathcal{X} $	the size (cardinality) of the set \mathcal{X}
$\text{Bern}(\lambda)$	Bernoulli (binary) random variable taking 1 with probability λ and 0 with probability $1 - \lambda$
$H_2(\lambda)$	entropy of $\text{Bern}(\lambda)$
$:=$	definition
Italic bold letters	Deterministic matrix \mathbf{M} / vector \mathbf{v}
Non-italic capital (bold) letters	Random variables X / Random vectors \mathbf{X}
$P(\cdot)$	Probability measure
$\mathbb{E}\{X\}$	Mean of the random variable X
$I(X; Y)$	Mutual information between X and Y
$\delta[\cdot]$	Kronecker delta function
$\mathbf{1}\{\cdot\}$	indicator function
$\log(x)$	Base-2 logarithm of x
\mathbf{I}	Identity matrix
$\ \mathbf{v}\ $	Euclidean (\mathcal{L}_2) norm of \mathbf{v}

Lecture 1: Information Measures

Lecturer: S. Yang

Disclaimer: These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.

The goal is to introduce the basic measures of information on which we rely throughout the course.

Probability measure and preliminaries

In this course, we consider a probability space (Ω, \mathcal{H}, P) where Ω is the **sample space**, \mathcal{H} is the σ -**algebra**, and P is the probability measure.

Let (E, \mathcal{E}) be a measurable space. A random variable X is a mapping $\Omega \rightarrow E$ that is measurable relative to \mathcal{H} and \mathcal{E} . In particular, if E is countable, then we call X a **discrete random variable**. For any $A \in \mathcal{E}$, we can define $\mu(A) := P(X(\omega) \in A)$, which is a probability measure on (E, \mathcal{E}) . For discrete random variables, we call $P_X(x) := P(X(\omega) = x)$, $x \in E$, the **probability mass function (pmf)**.

Let μ and ν be two measures on a measurable space (E, \mathcal{E}) , then ν is said to be **absolutely continuous** with respect to μ , denoted by $\nu \ll \mu$, if, for every set $A \in \mathcal{E}$, $\mu(A) = 0 \Rightarrow \nu(A) = 0$. If $\nu \ll \mu$, then there exists a Radon-Nikodym derivative of ν with respect to μ , often denoted by $\frac{d\nu}{d\mu}$, such that

$$\int_A \nu(dx) = \int_A \mu(dx) \frac{d\nu}{d\mu}(x), \quad \forall A \in \mathcal{E}.$$

Note that the Radon-Nikodym derivative is positive and measurable, i.e., in \mathcal{E}_+ . If $\nu \ll \mu \ll \lambda$, we have $\frac{d\nu}{d\mu} = \frac{d\nu}{d\lambda} \frac{d\mu}{d\lambda}$. If $P \ll Q$ are two probability measures defined on the same space (E, \mathcal{E}) , and f is a P -measurable function, then we have the change of measure $\mathbb{E}_P(f(X)) = \mathbb{E}_Q(f(X) \frac{dP}{dQ}(X))$.

Consider the case where E of the random variable X is the Euclidean space. If the probability measure μ is absolutely continuous with respect to the Lebesgue measure, then $p_X(x) := \frac{d\mu}{d\lambda}(x)$ is called the **probability density function (pdf)**. We call the random variable X a **continuous random variable**.

The mapping $(x, B) \mapsto K(x, B)$, $x \in E$ and $B \in \mathcal{F}$, is a **transition kernel** from (E, \mathcal{E}) to (F, \mathcal{F}) . In particular, we consider **probability transition kernel** such that $K(x, \mathcal{F}) = 1$ for all $x \in E$. If μ is a probability measure in E , then $\pi f = \int_E \mu(dx) \int_F K(x, dy) f(x, y)$ defines the unique probability measure satisfying $\pi(A \times B) = \int_A \mu(dx) K(x, B)$ for all $A \in \mathcal{E}, B \in \mathcal{F}$. Conversely, under some regularity conditions, for every probability measure on the product space $(E \times F, \mathcal{E} \otimes \mathcal{F})$, there exist a probability measure μ on E and a transition probability kernel K from (E, \mathcal{E}) to (F, \mathcal{F}) such that $\int_{E \times F} \pi(dx \times dy) f(x, y) = \int_E \mu(dx) \int_F K(x, dy) f(x, y)$, also known as “disintegration”. Throughout the course, we assume that such regularity conditions are met and ignore all measurability issues whenever possible. In most cases, we use $P_{Y|X}$ to denote the transition probability kernel such that $P_{Y|X=x}(dy) = K(x, dy)$. We use $P_X P_{Y|X}$ to denote the measure π on the product space such that $\int_{E \times F} \pi(dx \times dy) f(x, y) = \int_E P_X(dx) \int_F P_{Y|X=x}(dy) f(x, y)$. Both notations $P_{Y|X}(dy|x)$ and $P_{Y|X=x}(dy)$ are equivalent and can be used interchangeably. In particular, in the discrete case, $P_{Y|X}(y|x) = P_{Y|X=x}(y)$.

We use $P_{Y|X} \circ P_X$ to refer to the probability measure generated by the measure P_X and the transition kernel $P_{Y|X}$:

$$(P_{Y|X} \circ P_X)(A) = \int_E P_X(dx) P_{Y|X=x}(A).$$

It is the **marginalization** of the joint measure $P_{Y|X}P_X$. It can be regarded as the mixture of different distributions $P_{Y|X=x}$ according to the measure P_X . In the discrete case, the transition probability kernel is a matrix and the pmf's are column vectors, and \circ be done as matrix multiplications. In most cases, we use P and Q to denote probability measures and p and q as the corresponding density function, i.e., pmf in the discrete case (w.r.t. the counting measure) and pdf in the continuous case (w.r.t. the Lebesgue measure). Finally, we remove the subscript of the pmf/pdf whenever ambiguity is not likely.

We will use the terms *distribution* and *probability measure* interchangeably. Unless the context makes it obvious, the underlying probability distribution will always be specified. Given a joint distribution P_{XY} , one can derive the marginals P_X and P_Y , as well as the conditional distributions (transition kernels) $P_{Y|X}$ and $P_{X|Y}$, such that

$$P_{XY} = P_X P_{Y|X} = P_Y P_{X|Y}.$$

In this case, the notation $P_X, P_Y, P_{Y|X}, P_{X|Y}$ should always be understood as quantities induced by the given joint law P_{XY} . In more general situations, when several distributions or transition kernels are involved, we must be explicit about which objects are given and how others are constructed. For example, suppose we are given a conditional distribution $P_{Y|X}$ but not a joint law. Together with two different input distributions P_X and Q_X , we can form two distinct joint distributions, $P_{XY} := P_X P_{Y|X}$, $Q_{XY} := Q_X P_{Y|X}$. In this case, P_{XY} and Q_{XY} are simply notations with explicit definition from the original distributions and transition kernels.

A real function f is called **convex** in a set \mathcal{X} if for all $\lambda \in [0, 1]$ and all $x_1, x_2 \in \mathcal{X}$, we have

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2),$$

which is quite easy to visualize. A real function is called **concave** if $-f$ is convex, or, equivalently, the above inequality changes direction. For example, $x \mapsto \log(x)$ is concave in $(0, \infty)$, $x \mapsto x \log(x)$ is convex in $(0, \infty)$.

One of the most important inequalities that we use in information theory is the so-called **Jensen's inequality**: Let $X \in \mathcal{X}$ and f is convex in \mathcal{X} , then $\mathbb{E}f(X) \geq f(\mathbb{E}X)$. For concave functions, we have $\mathbb{E}f(X) \leq f(\mathbb{E}X)$.

1.1 Entropy

Now, we introduce the first information measure. The **entropy** $H(X)$ of a discrete random variable $X \sim P_X$ is defined as

$$H(X) \equiv H(P_X) := \mathbb{E}_{P_X} \log \frac{1}{p(X)} = \sum_{x \in \mathcal{X}} p(x) \log \frac{1}{p(x)}.$$

Sometimes, we use the notation $H(P_X)$ to emphasize that entropy is a functional of the pmf. Intuitively, entropy measures the *uncertainty* (or amount of *information*) of a random variable.

Similarly, we define the **joint entropy** with the joint distribution $P_{X_1 \dots X_n}$ (or, in short, P_{X^n}).

$$H(X_1, \dots, X_n) := \mathbb{E}_{P_{X^n}} \log \frac{1}{p(X_1, \dots, X_n)} = \sum_{x_1 \in \mathcal{X}_1} \dots \sum_{x_n \in \mathcal{X}_n} p(x_1, \dots, x_n) \log \frac{1}{p(x_1, \dots, x_n)}.$$

We also define the **conditional entropy** as

$$H(X|Y) = H(P_{X|Y} | P_Y) := \mathbb{E}_{y \sim P_Y} H(P_{X|Y=y}) = \mathbb{E}_{P_Y P_{X|Y}} \log \frac{1}{p(X|Y)}.$$

Here Y need not be discrete. If Y is also discrete, then

$$H(X|Y) := \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} p(x, y) \log \frac{1}{p(x|y)}.$$

For discrete $X^n \sim P_{X^n}$, we have the following **chain rule**:

$$H(X^n) = \sum_{i=1}^n H(X_i | X^{i-1})$$

This can be shown with the chain rule of probability $p(x^n) = \prod_i p(x_i | x^{i-1})$:

$$\begin{aligned} H(X^n) &= \mathbb{E}_{P_{X^n}} \left[\log \frac{1}{p(X_1, \dots, X_n)} \right] \\ &= \sum_{i=1}^n \mathbb{E}_{P_{X^n}} \left[\log \frac{1}{p(X_i | X^{i-1})} \right] \\ &= \sum_{i=1}^n \mathbb{E}_{P_{X^i}} \left[\log \frac{1}{p(X_i | X^{i-1})} \right] \\ &= \sum_{i=1}^n H(X_i | X^{i-1}) \end{aligned}$$

Similarly, the conditional version also holds.

$$H(X^n | Y) = \sum_{i=1}^n H(X_i | X^{i-1}, Y)$$

1.2 Cross-entropy

Let us consider two probability measures P and Q defined on the same measurable space. Define the **cross-entropy** of P with respect to Q as¹

$$H(P \| Q) := \begin{cases} \mathbb{E}_P \log \frac{1}{q(X)}, & P \ll Q \\ +\infty, & P \not\ll Q \end{cases}$$

And the **conditional cross-entropy**

$$H(P_{X|Y} \| Q_{X|Y} | P_Y) = \mathbb{E}_{y \sim P_Y} [H(P_{X|Y=y} \| Q_{X|Y=y})]$$

When $P_{X|Y=y} \ll Q_{X|Y=y}$, P_Y -a.s., we have

$$H(P_{X|Y} \| Q_{X|Y} | P_Y) = \mathbb{E}_{P_Y P_{X|Y}} \left[\log \frac{1}{q(X|Y)} \right]$$

As for entropy, the chain rule also holds for cross-entropy:

$$H(P_{X^n} \| Q_{X^n}) = \sum_{i=1}^n H(P_{X_i | X^{i-1}} \| Q_{X_i | X^{i-1}} | P_{X^{i-1}})$$

¹Here, we use a non-standard notation $H(P \| Q)$ instead of $H(P, Q)$ to avoid confusion with the joint entropy $H(X, Y)$.

$$\begin{aligned}
H(P_{X^n} \| Q_{X^n}) &= \mathbb{E}_{P_{X^n}} \left[\log \frac{1}{q(X_1, \dots, X_n)} \right] \\
&= \sum_{i=1}^n \mathbb{E}_{P_{X^n}} \left[\log \frac{1}{q(X_i | X^{i-1})} \right] \\
&= \sum_{i=1}^n \mathbb{E}_{P_{X^i}} \left[\log \frac{1}{q(X_i | X^{i-1})} \right] \\
&= \sum_{i=1}^n H(P_{X_i | X^{i-1}} \| Q_{X_i | X^{i-1}} | P_{X^{i-1}})
\end{aligned}$$

And the conditional version also holds.

$$H(P_{X^n | Y} \| Q_{X^n | Y} | P_Y) = \sum_{i=1}^n H(P_{X_i | X^{i-1} Y} \| Q_{X_i | X^{i-1} Y} | P_{X^{i-1} Y})$$

1.3 Differential entropy, differential cross-entropy

There is an equivalent definition of the entropy for continuous random variables. We define the **differential entropy** of a continuous random variable $X \sim P_X$ with pdf p as

$$h(X) \equiv h(P_X) := \mathbb{E}_{P_X} \left[\log \frac{1}{p(X)} \right] = \int dx p(x) \log \frac{1}{p(x)}.$$

Similarly, we define the conditional differential entropy

$$h(X | Y) \equiv h(P_{X|Y} | P_Y) := \mathbb{E}_{Y \sim P_Y} [h(P_{X|Y=y})] = \mathbb{E}_{P_Y P_{X|Y}} \left[\log \frac{1}{p(X|Y)} \right].$$

Finally, if X^n has a density, we also define the **joint differential entropy**

$$h(X^n) := \mathbb{E} \left[\log \frac{1}{p(X^n)} \right].$$

The differential cross-entropy $h(P \| Q)$ and the conditional version $h(P_{X|Y} \| Q_{X|Y} | P_Y)$ are defined as the cross-entropy but with the pdf's.

The chain rule holds as for entropy.

1.4 Divergence

Let us now introduce the divergence (aka. Kullback-Leibler divergence, KL divergence, or relative entropy). Consider two probability measures P and Q defined on the same measurable space. The **divergence** of P from Q , denoted by $D(P \| Q)$, is defined as

$$D(P \| Q) := \begin{cases} \mathbb{E}_P \left[\log \frac{dP}{dQ}(X) \right] = \mathbb{E}_Q \left[\frac{dP}{dQ}(X) \log \frac{dP}{dQ}(X) \right], & P \ll Q \\ +\infty, & P \not\ll Q. \end{cases}$$

If P and Q are discrete distributions defined on the same set \mathcal{X} , then the divergence becomes

$$D(P \| Q) := \begin{cases} \mathbb{E}_P \left[\log \frac{p(X)}{q(X)} \right] = \sum_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)}, & P \ll Q \\ +\infty, & P \not\ll Q \end{cases}$$

where the Radon-Nikodym derivative is simply the ratio of two pmf's.

If P and Q are distributions of continuous random variables defined on the same set \mathcal{X} , then the divergence becomes

$$D(P\|Q) := \begin{cases} \mathbb{E}_P \left[\log \frac{p(x)}{q(x)} \right] = \int_{\mathcal{X}} dx p(x) \log \frac{p(x)}{q(x)}, & P \ll Q \\ +\infty, & P \not\ll Q \end{cases}$$

where the Radon-Nikodym derivative is the ratio of two pdf's.

As for cross-entropy, one can also define the **conditional divergence** for an arbitrary probability measure P_Y and arbitrary probability transition kernels $P_{X|Y}$ and $Q_{X|Y}$,

$$D(P_{X|Y} \| Q_{X|Y} \mid P_Y) = \mathbb{E}_{y \sim P_Y} [D(P_{X|Y=y} \| Q_{X|Y=y})]$$

When $P_{X|Y=y} \ll Q_{X|Y=y}$, P_Y -a.s., we have

$$D(P_{X|Y} \| Q_{X|Y} \mid P_Y) = \mathbb{E}_{P_Y P_{X|Y}} \left[\log \frac{p(X|Y)}{q(X|Y)} \right]$$

Exactly as for cross-entropy, the chain rule of divergence holds:

$$D(P_{X^n} \| Q_{X^n}) = \sum_{i=1}^n D(P_{X_i|X^{i-1}} \| Q_{X_i|X^{i-1}} \mid P_{X^{i-1}})$$

and the conditional version

$$D(P_{X^n|Y} \| Q_{X^n|Y} \mid P_Y) = \sum_{i=1}^n D(P_{X_i|X^{i-1}Y} \| Q_{X_i|X^{i-1}Y} \mid P_{X^{i-1}Y})$$

For discrete distributions, the following relation between the entropy, cross-entropy, and divergence is straight-forward

$$H(P_{X|Y} \| Q_{X|Y} \mid P_Y) = H(P_{X|Y} \mid P_Y) + D(P_{X|Y} \| Q_{X|Y} \mid P_Y).$$

And in particular,

$$D(P\|Q) = H(P\|Q) - H(P).$$

The same holds for the continuous counterpart, by replacing entropy/cross-entropy by differential entropy and differential cross-entropy.

1.5 Mutual information

Let P_{XY} be the **joint distribution** of (X, Y) (e.g. probability measure of the product space $(E \times F, \mathcal{E} \otimes \mathcal{F})$). Similarly, let P_X and P_Y be the **marginal distributions** of X and Y , respectively. Further, let $P_{X|Y}$ and $P_{Y|X}$ be the transition probability kernels such that $P_{XY} = P_X P_{Y|X} = P_Y P_{X|Y}$.²

Then, **mutual information** measures the dependence between X and Y :

$$I(X; Y) := D(P_{XY} \| P_X P_Y) = D(P_{X|Y} \| P_X \mid P_Y) = D(P_{Y|X} \| P_Y \mid P_X),$$

where the last two equalities can be proved with the chain rule of divergence. It compares the joint distribution to the one where X and Y are independent, or the conditional distributions $P_{Y|X}$ and $P_{X|Y}$ to the marginals P_Y and P_X . Sometimes mutual information is also denoted by $I(P_X, P_{Y|X})$, $I(P_Y, P_{X|Y})$, and $I(P_{XY})$

$$I(X; Y) \equiv I(P_{XY}) \equiv I(P_X, P_{Y|X}) \equiv I(P_Y, P_{X|Y}).$$

²We assume that both kernels exists, which is guaranteed when both (E, \mathcal{E}) and (F, \mathcal{F}) are *standard spaces*.

Indeed, mutual information is a functional of the joint distribution P_{XY} .

In particular, in the discrete case

$$I(X; Y) = \mathbb{E}_{P_{XY}} \log \frac{p_{XY}(X, Y)}{p_X(X)p_Y(Y)}.$$

It follows that

$$I(X; Y) = H(X) + H(Y) - H(X, Y) \quad (1.1)$$

$$= H(X) - H(X|Y) \quad (1.2)$$

$$= H(Y) - H(Y|X) \quad (1.3)$$

The relationship between entropy and mutual information is best visualized with the Venn diagram below.

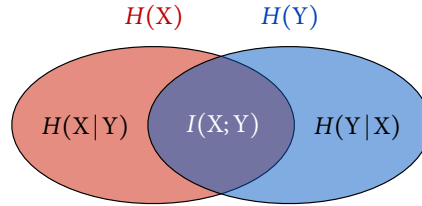


Figure 1.1: The Venn diagram.

In the continuous case, the definition is the same with the pdf

$$I(X; Y) = \mathbb{E}_{P_{XY}} \log \frac{p_{XY}(X, Y)}{p_X(X)p_Y(Y)}.$$

It follows that

$$I(X; Y) = h(X) + h(Y) - h(X, Y) \quad (1.4)$$

$$= h(X) - h(X|Y) \quad (1.5)$$

$$= h(Y) - h(Y|X), \quad (1.6)$$

similar to the discrete case.

If X is discrete and Y continuous, then

$$I(X; Y) = D(P_{XY} \| P_X P_Y) \quad (1.7)$$

$$= D(P_{Y|X} \| P_Y | P_X) \quad (1.8)$$

$$= h(Y) - h(Y|X). \quad (1.9)$$

We also have

$$I(X; Y) = H(X) - H(X|Y).$$

Note that in this case, although the conditional (differential) entropy exist, neither joint entropy nor joint differential entropy exists for (X, Y) .

Let P_{XYZ} be some joint distribution of (X, Y, Z) . Then, we can define the **conditional mutual information** between X and Y given Z .

$$I(X; Y | Z) \equiv I(P_{XY|Z} | P_Z) := D(P_{XY|Z} \| P_{X|Z} P_{Y|Z} | P_Z) = D(P_{X|YZ} \| P_{X|Z} | P_{YZ}) = D(P_{Y|XZ} \| P_{Y|Z} | P_{XZ}),$$

where we replace the divergence in the definition of mutual information by the conditional divergence given Y .

The chain rule of mutual information is

$$I(X; Y^n) = \sum_{i=1}^n I(X; Y_i | Y^{i-1})$$

Indeed, this can be proved from the (conditional) chain rule of divergence

$$\begin{aligned} I(X; Y^n) &= D(P_{Y^n|X} \| P_{Y^n} | P_X) \\ &= \sum_{i=1}^n D(P_{Y_i|XY^{i-1}} \| P_{Y_i|Y^{i-1}} | P_{XY^{i-1}}) \\ &= \sum_{i=1}^n I(X; Y_i | Y^{i-1}) \end{aligned}$$

The conditional version follows in the same way.

$$I(X; Y^n | Z) = \sum_{i=1}^n I(X; Y_i | Y^{i-1}Z)$$

1.6 Some properties of information measures

- General chain rule: Writing the chain rules in the same notational convention, we have

$$\begin{aligned} H(P_{X^n}) &= \sum_{i=1}^n H(P_{X_i|X^{i-1}} | P_{X^{i-1}}), \quad h(P_{X^n}) = \sum_{i=1}^n h(P_{X_i|X^{i-1}} | P_{X^{i-1}}) \\ H(P_{X^n} \| Q_{X^n}) &= \sum_{i=1}^n H(P_{X_i|X^{i-1}} \| Q_{X_i|X^{i-1}} | P_{X^{i-1}}), \quad h(P_{X^n} \| Q_{X^n}) = \sum_{i=1}^n h(P_{X_i|X^{i-1}} \| Q_{X_i|X^{i-1}} | P_{X^{i-1}}) \\ D(P_{X^n} \| Q_{X^n}) &= \sum_{i=1}^n D(P_{X_i|X^{i-1}} \| Q_{X_i|X^{i-1}} | P_{X^{i-1}}) \\ I(P_{X, Y^n}) &= \sum_{i=1}^n I(P_{X, Y_i|Y^{i-1}} | P_{Y^{i-1}}) \end{aligned}$$

- Positivity

$$\begin{aligned} H(P) &\geq 0, \quad h(P) \not\leq 0 \\ H(P \| Q) &\geq 0, \quad h(P \| Q) \not\leq 0 \\ D(P \| Q) &\geq 0 \implies H(P) \leq H(P \| Q), \quad h(P) \leq h(P \| Q) \\ I(P_X, P_{Y|X}) &\geq 0 \end{aligned}$$

Proof. For entropy, since probability is upper bounded by 1, entropy is nonnegative. Entropy is 0 if and only if the random variable is deterministic. The positivity does not hold for differential entropy. The same arguments apply for cross-entropy and differential cross-entropy.

For divergence, if $P \not\ll Q$, then $D(P \| Q) = +\infty > 0$. We assume therefore $P \ll Q$. Then, let us write $D(P \| Q) = \mathbb{E}_Q \left(f \left(\frac{dP}{dQ} \right) \right)$ where $f(x) := x \log x$. Finally, since $f(x)$ is strictly convex (check), we have $D(P \| Q) \geq f(\mathbb{E}_Q \frac{dP}{dQ}) = f(1) = 0$, where we applied Jensen's inequality on f . The equality holds if and only if $\frac{dP}{dQ}$ is constant (Q -almost everywhere), implying that $P = Q$.

The positivity of mutual information is from that of divergence. It is 0 if and only if $P_{XY} = P_X P_Y$, i.e., X and Y are independent.

□

- Conditioning

- Conditioning reduces (differential) entropy

$$H(X) \geq H(X|Y), \quad h(X) \geq h(X|Y)$$

- Conditioning increases divergence

$$D(P_{X|Y} \| Q_{X|Y} | P_Y) \geq D(\tilde{P}_X \| \tilde{Q}_X)$$

where \tilde{P}_X and \tilde{Q}_X are the marginals of $P_{X|Y}P_Y$ and $Q_{X|Y}P_Y$ respectively, i.e., $\tilde{P}_X = \mathbb{E}_{y \sim P_Y}[P_{X|Y=y}]$ and $\tilde{Q}_X = \mathbb{E}_{y \sim P_Y}[Q_{X|Y=y}]$.

Proof. Conditioning reduces entropy is from the positivity of mutual information, i.e., $H(X) - H(X|Y) = I(X;Y) \geq 0$. Similarly for differential entropy. For divergence, we have $D(P_{X|Y} \| Q_{X|Y} | P_Y) = D(P_{X|Y}P_Y \| Q_{X|Y}P_Y) = D(\tilde{P}_X \| \tilde{Q}_X) + D(\tilde{P}_{Y|X} \| \tilde{Q}_{Y|X} | \tilde{P}_X) \geq D(\tilde{P}_X \| \tilde{Q}_X)$, where we use the decompositions $P_{X|Y}P_Y = \tilde{P}_X\tilde{P}_{Y|X}$ and $Q_{X|Y}P_Y = \tilde{Q}_X\tilde{Q}_{Y|X}$. \square

- Convexity/Concavity

- $P \mapsto H(P)$, $P \mapsto h(P)$ are both concave
- $(P, Q) \mapsto D(P \| Q)$ is convex
- $P_X \mapsto I(P_X, P_{Y|X})$ is concave
- $P_{Y|X} \mapsto I(P_X, P_{Y|X})$ is convex

Proof. Fix $\lambda \in [0, 1]$. Let $S \sim P_S := \text{Bern}(\lambda)$, i.e., $p_S(0) = \lambda$ and $p_S(1) = 1 - \lambda$.

Let $P_{X|S=0} := P_0$ and $P_{X|S=1} := P_1$. We have $P_X = (1 - \lambda)P_0 + \lambda P_1$. $\lambda H(P_1) + (1 - \lambda)H(P_0) = H(X|S) \leq H(X) = H(P_X)$, proving the concavity of $P \mapsto H(P)$ using conditioning reduces entropy.

For divergence, let $P_{X|S=k} := P_k$ and $Q_{X|S=k} := Q_k$ for $k = 0, 1$. Then, $\lambda D(P_1 \| Q_1) + (1 - \lambda)D(P_0 \| Q_0) = (P_{X|S} \| Q_{X|S} | P_S) \geq D(P_X \| Q_X)$, proving the convexity of divergence using conditioning increases divergence.

For the concavity of mutual information, for the given P_0, P_1 , and $P_{Y|X}$, let us set $P_{X|S=0} := P_0$ and $P_{X|S=1} := P_1$, and let $P_{Y|XS} = P_{Y|X}$, i.e., $P_{Y|X,S=0} = P_{Y|X,S=1} = P_{Y|X}$. Hence, the joint distribution is $P_{SXY} = P_S P_{X|S} P_{Y|X}$. The conditional mutual information

$$\begin{aligned} I(X; Y | S) &= D(P_{X|Y|S} \| P_{X|S} P_{Y|S} | P_S) \\ &= \lambda D(P_{Y|X,S=0} \| P_{Y|S=0} | P_{X|S=0}) + (1 - \lambda) D(P_{Y|X,S=1} \| P_{Y|S=1} | P_{X|S=1}) \\ &= \lambda D(P_{Y|X} \| P_{Y|S=0} | P_0) + (1 - \lambda) D(P_{Y|X} \| P_{Y|S=1} | P_1) \\ &= \lambda I(P_0, P_{Y|X}) + (1 - \lambda) I(P_1, P_{Y|X}). \end{aligned}$$

On the other hand,

$$I(X; Y) = I(P_X, P_{Y|X}) = I(\lambda P_0 + (1 - \lambda)P_1, P_{Y|X})$$

To finish the proof, we write

$$\begin{aligned} I(X; Y) &= I(X; Y) + I(S; Y | X) \\ &= I(X, S; Y) \\ &= I(S; Y) + I(X; Y | S) \\ &\geq I(X; Y | S) \end{aligned}$$

where the first equality holds since $I(S; Y | X) = 0$ due to the Markov chain $S \rightarrow X \rightarrow Y$. Indeed,

$$\begin{aligned} I(S; Y | X) &= D(P_{Y|XS}; P_{Y|X} | P_{XS}) \\ &= D(P_{Y|X}; P_{Y|X} | P_{XS}) \\ &= 0 \end{aligned}$$

Finally, for the convexity of mutual information, we need to prove that given P_X $I(P_X, \lambda P_{Y|X}^1 + (1 - \lambda)P_{Y|X}^0) \leq \lambda I(P_X, P_{Y|X}^1) + (1 - \lambda)I(P_X, P_{Y|X}^0)$ for any kernels $P_{Y|X}^0$ and $P_{Y|X}^1$ and $\lambda \in [0, 1]$. We can prove it in two different ways. First, we can apply the convexity of divergence. Indeed,

$$\begin{aligned} I(P_X, \lambda P_{Y|X}^1 + (1 - \lambda)P_{Y|X}^0) &= D(\lambda P_X P_{Y|X}^1 + (1 - \lambda)P_X P_{Y|X}^0 \| \lambda P_X P_Y^1 + (1 - \lambda)P_X P_Y^0) \\ &\geq \lambda D(P_X P_{Y|X}^1 \| P_X P_Y^1) + (1 - \lambda)D(P_X P_{Y|X}^0 \| P_X P_Y^0) \\ &= \lambda I(P_X, P_{Y|X}^1) + (1 - \lambda)I(P_X, P_{Y|X}^0) \end{aligned}$$

where P_Y^0 and P_Y^1 are the marginals of $P_X P_{Y|X}^0$ and $P_X P_{Y|X}^1$, respectively.

The second way is to introduce the same S as before, let $P_{Y|X,S=0} = P_{Y|X}^0$ and $P_{Y|X,S=1} = P_{Y|X}^1$, so that $P_{SXY} = P_S P_X P_{Y|XS}$. Unlike the previous cases, here X and S are independent. It can be verified that $P_{Y|X} = \lambda P_{Y|X}^0 + (1 - \lambda)P_{Y|X}^1$. Therefore, we have $I(P_X, \lambda P_{Y|X}^1 + (1 - \lambda)P_{Y|X}^0) = I(X; Y)$. We also have $\lambda I(P_X, P_{Y|X}^1) + (1 - \lambda)I(P_X, P_{Y|X}^0) = I(X; Y | S)$. Therefore, it is enough to prove that $I(X; Y | S) \geq I(X; Y)$. To that end, apply the independence so that $I(X; S) = 0$, and thus

$$\begin{aligned} I(X; Y | S) &= I(X; Y | S) + I(X; S) \\ &= I(X; Y, S) \\ &= I(X; Y) + I(S; Y | X) \\ &\geq I(X; Y). \end{aligned}$$

□

Exercises³

- Entropy of functions [CT 2.2]. Let X be a random variable taking on a finite number of values. What is the (general) inequality relationship of $H(X)$ and $H(Y)$ if
 - $Y = 2^X$?
 - $Y = \cos(X)$?
- Conditional mutual information vs. unconditional mutual information [CT 2.6]. Give examples of joint random variables X , Y , and Z such that
 - $I(X; Y | Z) < I(X; Y)$.
 - $I(X; Y | Z) > I(X; Y)$.

- Data processing [CT 2.15]. Let $X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow \dots \rightarrow X_n$ form a Markov chain in this order; that is, let

$$p(x_1, x_2, \dots, x_n) = p(x_1)p(x_2 | x_1) \cdots p(x_n | x_{n-1}).$$

Reduce $I(X_1; X_2, \dots, X_n)$ to its simplest form.

- Infinite entropy. [CT 2.19] This problem shows that the entropy of a discrete random variable can be infinite. Let $A = \sum_{n=2}^{\infty} (n \log^2 n)^{-1}$. [It is easy to show that A is finite by bounding the infinite sum by the integral of $(x \log^2 x)^{-1}$.] Show that the integer-valued random variable X defined by $P(X = n) = (An \log^2 n)^{-1}$ for $n = 2, 3, \dots$ has $H(X) = +\infty$.
- Inequalities [CT 2.29]. Let X , Y , and Z be joint random variables. Prove the following inequalities and find conditions for equality.
 - $H(X, Y | Z) \geq H(X | Z)$.
 - $I(X, Y; Z) \geq I(X; Z)$.
 - $H(X, Y, Z) - H(X, Y) \leq H(X, Z) - H(X)$.
 - $I(X; Z | Y) \geq I(Z; Y | X) - I(Z; Y) + I(X; Z)$.

- Convexity/Concavity of mutual information.

- Let $(S, X, Y) \sim P_{SXY} = P_S P_{X|S} P_{Y|X}$, i.e., $S \rightarrow X \rightarrow Y$ forms a Markov chain. Show that

$$I(X; Y) \geq I(X; Y | S).$$

Use the above inequality to show that mutual information is concave in P_X for a fixed $P_{Y|X}$.

- Let $(S, X, Y) \sim P_{SXY} = P_S P_X P_{Y|X, S}$. Show that

$$I(X; Y) \leq I(X; Y | S).$$

Use the above inequality to show that mutual information is convex in $P_{Y|X}$ for a fixed P_X .

- Maximum entropy. [CT 2.30] Find the probability mass function $P(x)$ that maximizes the entropy $H(X)$ of a nonnegative integer-valued random variable X subject to the constraint

$$E(X) = \sum_{n=0}^{\infty} nP(n) = A$$

for a fixed value $A > 0$. Evaluate this maximum $H(X)$.

- Relative entropy is not symmetric. [CT 2.35] Let the random variable X have three possible outcomes $\{a, b, c\}$. Consider two distributions on this random variable:

Symbol	$P(x)$	$Q(x)$
a	$\frac{1}{2}$	$\frac{1}{3}$
b	$\frac{1}{4}$	$\frac{1}{3}$
c	$\frac{1}{4}$	$\frac{1}{3}$

³The citations "CT", "CK" refer to Cover-Thomas, Csiszár-Körner, respectively.

9. Entropy and pairwise independence. [CT 2.39] Let X, Y, Z be three binary Bernoulli(1/2) random variables that are pairwise independent; that is, $I(X; Y) = I(X; Z) = I(Y; Z) = 0$.
- Under this constraint, what is the minimum value for constraint, $H(X, Y, Z)$?
 - Give an example achieving this minimum.
10. Mutual information of heads and tails [CT 2.43]
- Consider a fair coin flip. What is the mutual information between the top and bottom sides of the coin?
 - A six-sided fair die is rolled. What is the mutual information between the top side and the front face (the side most facing you)?
11. Finite entropy. [CT 2.45] Show that for a discrete random variable $X \in \{1, 2, \dots\}$, if $E \log X < \infty$, then $H(X) < \infty$.
12. Sequence length. [CT 2.48] How much information does the length of a sequence give about the content of a sequence? Suppose that we consider a Bernoulli(1/2) process $\{X_i\}$. Stop the process when the first 1 appears. Let N designate this stopping time. Thus, X^N is an element of the set of all finite-length binary sequences $\{0, 1\}^* = \{0, 1, 00, 01, 10, 11, 000, \dots\}$
- Find $I(N; X^N)$.
 - Find $H(X^N | N)$.
 - Find $H(X^N)$.
- Let's now consider a different stopping time. For this part, again assume that $X \sim \text{Bernoulli}(1/2)$ but stop at time $N = 6$, with probability $1/3$ and stop at time $N = 12$ with probability $2/3$. Let this stopping time be independent of the sequence X_1, X_2, \dots, X_{12} .
- Find $I(N; X^N)$. Find $H(X | N)$. Find $H(X^N)$.
13. Function of variables from a Markov chain. [CK 3.7] Is it true that if $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_n$, and f is an arbitrary function on the common range of the X_i 's, then $f(X_1) \rightarrow f(X_2) \rightarrow \dots \rightarrow f(X_n)$? Give a counter example.
14. Mutual information. [Gallager 2.8] Consider an ensemble of sequences of N binary digits, x_1, x_2, \dots, x_N . Each sequence containing an even number of 1's has probability 2^{-N+1} and each sequence with an odd number of 1's has probability zero. Find the mutual informations

$$I(X_1; X_2), I(X_2; X_3 | X_1), \dots, I(X_{N-1}; X_N | X_1, \dots, X_{N-2}).$$

Check your result for $N = 3$.

15. Memoryless source. Consider a sequence from the source contains independent symbols, i.e., $P_{X^n} = P_{X_1} \cdots P_{X_n}$, also denoted by $\prod_{i=1}^n P_{X_i}$. Show that

$$I(X^n; Y^n) \geq \sum_{i=1}^n I(X_i; Y_i),$$

for any $P_{Y^n|X^n}$, with equality if and only if $P_{X^n|Y^n} = \prod_i P_{X_i|Y_i}$ (P_{Y^n} -almost surely). *Hint: Apply chain rule on X^n , then use the independence between X_i and X^{i-1} .*

16. Memoryless channels without feedback. We say that a channel is memoryless *without feedback* if $P_{Y^n|X^n} = \prod_{i=1}^n P_{Y_i|X_i}$. Show that in this case we have the Markov chain $Y_i \rightarrow X_i \rightarrow (\{X_j, j \neq i\}, \{Y_j, j \neq i\})$. Show that

$$I(X^n; Y^n) \leq \sum_{i=1}^n I(X_i; Y_i),$$

with equality if and only if $P_{Y^n} = \prod_i P_{Y_i}$. *Hint: Apply the chain rule and the Markov chain.*

Quiz (unique correct answer)

- For a discrete random variable X taking on n possible values, which of the following statements is **TRUE** regarding its Shannon entropy $H(X)$?
 - $H(X)$ is maximized when X is deterministic.
 - $H(X)$ is minimized when X follows a uniform distribution.
 - $H(X)$ is always non-negative and less than or equal to $\log_2 n$.
 - $H(X)$ can be negative if X takes negative values.
 - $H(X)$ measures the variance of X .
- Which of the following expressions correctly represents the mutual information $I(X; Y)$ between two discrete random variables X and Y ?
 - $I(X; Y) = H(X, Y) - H(X) - H(Y)$
 - $I(X; Y) = H(X) + H(Y) - H(X, Y)$
 - $I(X; Y) = H(X|Y) - H(X)$
 - $I(X; Y) = H(X, Y) + H(X|Y)$
 - $I(X; Y) = H(X|Y) + H(Y|X)$
- Suppose X and Y are independent discrete random variables. Which of the following is **TRUE**?
 - $H(X|Y) = H(X)$
 - $H(X|Y) = 0$
 - $H(X, Y) = H(X)$
 - $I(X; Y) = H(X)$
 - $I(X; Y) = H(Y)$
- Which of the following is **TRUE** about the Kullback-Leibler divergence $D(P\|Q)$ between two discrete probability distributions P and Q ?
 - $D(P\|Q)$ is symmetric in P and Q .
 - $D(P\|Q) \geq 0$, and equals zero if and only if $P = Q$ almost everywhere.
 - $D(P\|Q)$ is always finite.
 - $D(P\|Q)$ measures the variance between P and Q .
 - $D(P\|Q)$ is the mutual information between $X \sim P$ and $Y \sim Q$.
- Which of the following information measures can be negative?
 - Shannon entropy $H(X)$
 - Mutual information $I(X; Y)$
 - Conditional entropy $H(X|Y)$
 - Differential entropy $h(X)$ of a continuous random variable X
 - Kullback-Leibler divergence $D(P\|Q)$
- For a continuous random variable X with probability density function $f(x)$, which of the following statements about the differential entropy $h(X)$ is **TRUE**?
 - $h(X)$ is always non-negative.
 - $h(X)$ is invariant under scaling of X .
 - $h(X)$ increases when X is scaled by a factor $a > 1$.

- D) $h(X)$ cannot be less than zero.
- E) $h(X)$ measures the variance of X .
7. The entropy $H(X)$ of a Bernoulli random variable X with parameter p (i.e., $P(X = 1) = p$) is given by:
- A) $H(X) = -p \log p$
- B) $H(X) = -p \log p - (1 - p) \log(1 - p)$
- C) $H(X) = p \log(1 - p) + (1 - p) \log p$
- D) $H(X) = -\log p$
- E) $H(X) = p$
8. Which of the following inequalities relates the conditional entropy $H(X|Y)$ and the entropy $H(X)$ of two discrete random variables X and Y ?
- A) $H(X|Y) \geq H(X)$
- B) $H(X|Y) \leq H(X)$
- C) $H(X|Y) = H(X)$
- D) $H(X|Y) = H(X) + H(Y)$
- E) $H(X|Y) = H(X, Y) - H(Y)$
9. The chain rule for entropy states that for discrete random variables X and Y :
- A) $H(X, Y) = H(X|Y) + H(Y)$
- B) $H(X, Y) = H(X) + H(Y)$
- C) $H(X, Y) = H(Y|X) - H(X)$
- D) $H(X, Y) = H(X|Y) - H(Y)$
- E) $H(X, Y) = H(X) - H(Y|X)$
10. Which of the following distributions maximizes the entropy among all continuous distributions with a given variance?
- A) Uniform distribution
- B) Exponential distribution
- C) Gaussian (Normal) distribution
- D) Laplace distribution
- E) Cauchy distribution
11. The conditional mutual information $I(X; Y|Z)$ can be expressed in terms of entropies as:
- A) $I(X; Y|Z) = H(X|Z) + H(Y|Z) - H(X, Y|Z)$
- B) $I(X; Y|Z) = H(X, Y, Z) - H(Z)$
- C) $I(X; Y|Z) = H(X|Y, Z) - H(X|Z)$
- D) $I(X; Y|Z) = H(X, Z) + H(Y, Z) - H(Z)$
- E) $I(X; Y|Z) = H(X, Y) - H(Z)$
12. For a discrete random variable X and function $f(X)$, which of the following is **TRUE** regarding the entropy $H(f(X))$?
- A) $H(f(X)) \geq H(X)$
- B) $H(f(X)) \leq H(X)$
- C) $H(f(X)) = H(X)$

- D) $H(f(X)) = 0$
 E) $H(f(X)) = H(X|f(X))$
13. The Data Processing Inequality states that for random variables forming a Markov chain $X \rightarrow Y \rightarrow Z$, which of the following is **TRUE**?
- A) $I(X; Z) \geq I(X; Y)$
 B) $I(X; Z) \leq I(X; Y)$
 C) $I(X; Y) = I(Y; Z)$
 D) $I(X; Z) = I(X; Y) + I(Y; Z)$
 E) $I(X; Z) \geq I(Y; Z)$
14. Which distribution maximizes the entropy for a discrete random variable X with a fixed mean over the support $\{1, 2, \dots, n\}$?
- A) Uniform distribution over $\{1, 2, \dots, n\}$
 B) Geometric distribution
 C) Binomial distribution
 D) Discrete exponential distribution
 E) Poisson distribution
15. Define the cross-entropy $H(P, Q) := \mathbb{E}_P \log \frac{1}{Q(X)}$. Which of the following is a property of the cross-entropy $H(P, Q)$ between two probability distributions P and Q ?
- A) $H(P, Q) = H(Q, P)$
 B) $H(P, Q) \geq H(P)$
 C) $H(P, Q) \leq H(P)$
 D) $H(P, Q) = H(P) + D(P\|Q)$
 E) $H(P, Q) = D(P\|Q)$
16. For two independent continuous random variables X and Y , the differential entropy of their sum $Z = X + Y$ satisfies:
- A) $h(Z) = h(X) + h(Y)$
 B) $h(Z) = h(X) - h(Y)$
 C) $h(Z) = h(X) + h(Y) + \log 2\pi e$
 D) $h(Z) \leq h(X) + h(Y)$
 E) $h(Z) \geq h(X) + h(Y)$
17. Which of the following statements about mutual information $I(X; Y)$ is **TRUE**?
- A) Mutual information $I(X; Y)$ is always less than or equal to zero.
 B) Mutual information $I(X; Y)$ is zero if and only if X and Y are independent.
 C) Mutual information $I(X; Y)$ is the same as conditional entropy $H(X|Y)$.
 D) Mutual information $I(X; Y)$ is maximized when X and Y are independent.
 E) Mutual information $I(X; Y)$ is always greater than the joint entropy $H(X, Y)$.
18. The conditional entropy $H(Y|X)$ can be expressed in terms of joint entropy $H(X, Y)$ and marginal entropy $H(X)$ as:
- A) $H(Y|X) = H(X, Y) - H(X)$
 B) $H(Y|X) = H(X) - H(X, Y)$

- C) $H(Y|X) = H(Y) - H(X)$
 - D) $H(Y|X) = H(X, Y) + H(X)$
 - E) $H(Y|X) = H(Y) + H(X, Y)$
19. For a continuous random variable X uniformly distributed over the interval $[a, b]$, the differential entropy $h(X)$ is:
- A) $h(X) = \log(b - a)$
 - B) $h(X) = \log(b + a)$
 - C) $h(X) = \frac{1}{2} \log(b - a)$
 - D) $h(X) = -\log(b - a)$
 - E) $h(X) = \log\left(\frac{b}{a}\right)$
20. The **Chain Rule** for mutual information states that for random variables X, Y, Z :
- A) $I(X; Y, Z) = I(X; Y) + I(X; Z)$
 - B) $I(X; Y, Z) = I(X; Y|Z) + I(X; Z)$
 - C) $I(X; Y, Z) = I(X; Y) + I(X; Z|Y)$
 - D) $I(X; Y, Z) = I(X; Y) - I(X; Z|Y)$
 - E) $I(X; Y, Z) = I(X; Y|Z) - I(X; Z|Y)$