

Chapter 2

Entropy, Relative Entropy and Mutual Information

This chapter introduces most of the basic definitions required for the subsequent development of the theory. It is irresistible to play with their relationships and interpretations, taking faith in their later utility. After defining entropy and mutual information, we establish chain rules, the non-negativity of mutual information, the data processing inequality, and finally investigate the extent to which the second law of thermodynamics holds for Markov processes.

The concept of information is too broad to be captured completely by a single definition. However, for any probability distribution, we define a quantity called the *entropy*, which has many properties that agree with the intuitive notion of what a measure of information should be. This notion is extended to define *mutual information*, which is a measure of the amount of information one random variable contains about another. Entropy then becomes the self-information of a random variable. Mutual information is a special case of a more general quantity called *relative entropy*, which is a measure of the distance between two probability distributions. All these quantities are closely related and share a number of simple properties. We derive some of these properties in this chapter.

In later chapters, we show how these quantities arise as natural answers to a number of questions in communication, statistics, complexity and gambling. That will be the ultimate test of the value of these definitions.

2.1 ENTROPY

We will first introduce the concept of entropy, which is a measure of uncertainty of a random variable. Let X be a discrete random variable

with alphabet \mathcal{X} and probability mass function $p(x) = \Pr\{X = x\}$, $x \in \mathcal{X}$. We denote the probability mass function by $p(x)$ rather than $p_X(x)$ for convenience. Thus, $p(x)$ and $p(y)$ refer to two different random variables, and are in fact different probability mass functions, $p_X(x)$ and $p_Y(y)$ respectively.

Definition: The *entropy* $H(X)$ of a discrete random variable X is defined by

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log p(x). \quad (2.1)$$

We also write $H(p)$ for the above quantity. The log is to the base 2 and entropy is expressed in bits. For example, the entropy of a fair coin toss is 1 bit. We will use the convention that $0 \log 0 = 0$, which is easily justified by continuity since $x \log x \rightarrow 0$ as $x \rightarrow 0$. Thus adding terms of zero probability does not change the entropy.

If the base of the logarithm is b , we will denote the entropy as $H_b(X)$. If the base of the logarithm is e , then the entropy is measured in *nats*. Unless otherwise specified, we will take all logarithms to base 2, and hence all the entropies will be measured in bits.

Note that entropy is a functional of the distribution of X . It does not depend on the actual values taken by the random variable X , but only on the probabilities.

We shall denote expectation by E . Thus if $X \sim p(x)$, then the expected value of the random variable $g(X)$ is written

$$E_p g(X) = \sum_{x \in \mathcal{X}} g(x) p(x), \quad (2.2)$$

or more simply as $Eg(X)$ when the probability mass function is understood from the context.

We shall take a peculiar interest in the eerily self-referential expectation of $g(X)$ under $p(x)$ when $g(X) = \log \frac{1}{p(X)}$.

Remark: The entropy of X can also be interpreted as the expected value of $\log \frac{1}{p(X)}$, where X is drawn according to probability mass function $p(x)$. Thus

$$H(X) = E_p \log \frac{1}{p(X)}. \quad (2.3)$$

This definition of entropy is related to the definition of entropy in thermodynamics; some of the connections will be explored later. It is possible to derive the definition of entropy axiomatically by defining certain properties that the entropy of a random variable must satisfy. This approach is illustrated in a problem at the end of the chapter. We

will not use the axiomatic approach to justify the definition of entropy; instead, we will show that it arises as the answer to a number of natural questions such as “What is the average length of the shortest description of the random variable?” First, we derive some immediate consequences of the definition.

Lemma 2.1.1: $H(X) \geq 0$.

Proof: $0 \leq p(x) \leq 1$ implies $\log(1/p(x)) \geq 0$. \square

Lemma 2.1.2: $H_b(X) = (\log_b a)H_a(X)$.

Proof: $\log_b p = \log_b a \log_a p$. \square

The second property of entropy enables us to change the base of the logarithm in the definition. Entropy can be changed from one base to another by multiplying by the appropriate factor.

Example 2.1.1: Let

$$X = \begin{cases} 1 & \text{with probability } p, \\ 0 & \text{with probability } 1 - p. \end{cases} \quad (2.4)$$

Then

$$H(X) = -p \log p - (1 - p) \log(1 - p) \stackrel{\text{def}}{=} H(p). \quad (2.5)$$

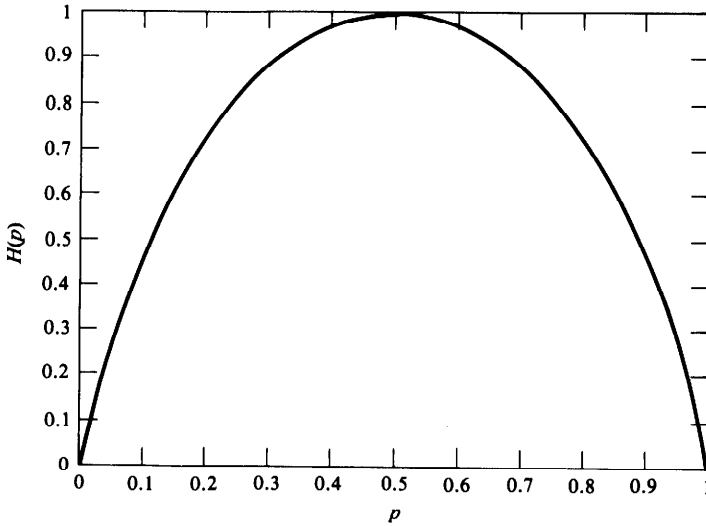
In particular, $H(X) = 1$ bit when $p = 1/2$. The graph of the function $H(p)$ is shown in Figure 2.1. The figure illustrates some of the basic properties of entropy—it is a concave function of the distribution and equals 0 when $p = 0$ or 1. This makes sense, because when $p = 0$ or 1, the variable is not random and there is no uncertainty. Similarly, the uncertainty is maximum when $p = \frac{1}{2}$, which also corresponds to the maximum value of the entropy.

Example 2.1.2: Let

$$X = \begin{cases} a & \text{with probability } 1/2, \\ b & \text{with probability } 1/4, \\ c & \text{with probability } 1/8, \\ d & \text{with probability } 1/8. \end{cases} \quad (2.6)$$

The entropy of X is

$$H(X) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{8} \log \frac{1}{8} - \frac{1}{8} \log \frac{1}{8} = \frac{7}{4} \text{ bits}. \quad (2.7)$$

Figure 2.1. $H(p)$ versus p .

Suppose we wish to determine the value of X with the minimum number of binary questions. An efficient first question is “Is $X = a$?” This splits the probability in half. If the answer to the first question is no, then the second question can be “Is $X = b$?” The third question can be “Is $X = c$?” The resulting expected number of binary questions required is 1.75. This turns out to be the minimum expected number of binary questions required to determine the value of X . In Chapter 5, we show that the minimum expected number of binary questions required to determine X lies between $H(X)$ and $H(X) + 1$.

2.2 JOINT ENTROPY AND CONDITIONAL ENTROPY

We have defined the entropy of a single random variable in the previous section. We now extend the definition to a pair of random variables. There is nothing really new in this definition because (X, Y) can be considered to be a single vector-valued random variable.

Definition: The *joint entropy* $H(X, Y)$ of a pair of discrete random variables (X, Y) with a joint distribution $p(x, y)$ is defined as

$$H(X, Y) = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x, y), \quad (2.8)$$

which can also be expressed as

$$H(X, Y) = -E \log p(X, Y). \quad (2.9)$$

We also define the conditional entropy of a random variable given another as the expected value of the entropies of the conditional distributions, averaged over the conditioning random variable.

Definition: If $(X, Y) \sim p(x, y)$, then the *conditional entropy* $H(Y|X)$ is defined as

$$H(Y|X) = \sum_{x \in \mathcal{X}} p(x) H(Y|X=x) \quad (2.10)$$

$$= - \sum_{x \in \mathcal{X}} p(x) \sum_{y \in \mathcal{Y}} p(y|x) \log p(y|x) \quad (2.11)$$

$$= - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(y|x) \quad (2.12)$$

$$= - E_{p(x, y)} \log p(Y|X). \quad (2.13)$$

The naturalness of the definition of joint entropy and conditional entropy is exhibited by the fact that the entropy of a pair of random variables is the entropy of one plus the conditional entropy of the other. This is proved in the following theorem.

Theorem 2.2.1 (*Chain rule*):

$$H(X, Y) = H(X) + H(Y|X). \quad (2.14)$$

Proof:

$$H(X, Y) = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x, y) \quad (2.15)$$

$$= - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x) p(y|x) \quad (2.16)$$

$$= - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x) - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(y|x) \quad (2.17)$$

$$= - \sum_{x \in \mathcal{X}} p(x) \log p(x) - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(y|x) \quad (2.18)$$

$$= H(X) + H(Y|X). \quad (2.19)$$

Equivalently, we can write

$$\log p(X, Y) = \log p(X) + \log p(Y|X) \quad (2.20)$$

and take the expectation of both sides of the equation to obtain the theorem. \square

Corollary:

$$H(X, Y|Z) = H(X|Z) + H(Y|X, Z). \quad (2.21)$$

Proof: The proof follows along the same lines as the theorem. \square

Example 2.2.1: Let (X, Y) have the following joint distribution:

Y \ X	1	2	3	4
1	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{32}$	$\frac{1}{32}$
2	$\frac{1}{16}$	$\frac{1}{8}$	$\frac{1}{32}$	$\frac{1}{32}$
3	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$
4	$\frac{1}{4}$	0	0	0

The marginal distribution of X is $(\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8})$ and the marginal distribution of Y is $(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4})$, and hence $H(X) = 7/4$ bits and $H(Y) = 2$ bits. Also,

$$\begin{aligned} H(X|Y) &= \sum_{i=1}^4 p(Y=i)H(X|Y=i) \\ &= \frac{1}{4} H\left(\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8}\right) + \frac{1}{4} H\left(\frac{1}{4}, \frac{1}{2}, \frac{1}{8}, \frac{1}{8}\right) \end{aligned} \quad (2.22)$$

$$+ \frac{1}{4} H\left(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}\right) + \frac{1}{4} H(1, 0, 0, 0) \quad (2.23)$$

$$= \frac{1}{4} \times \frac{7}{4} + \frac{1}{4} \times \frac{7}{4} + \frac{1}{4} \times 2 + \frac{1}{4} \times 0 \quad (2.24)$$

$$= \frac{11}{8} \text{ bits.} \quad (2.25)$$

Similarly $H(Y|X) = 13/8$ bits and $H(X, Y) = 27/8$ bits.

Remark: Note that $H(Y|X) \neq H(X|Y)$. However, $H(X) - H(X|Y) = H(Y) - H(Y|X)$, a property that we shall exploit later.

2.3 RELATIVE ENTROPY AND MUTUAL INFORMATION

The entropy of a random variable is a measure of the uncertainty of the random variable; it is a measure of the amount of information required on the average to describe the random variable. In this section, we introduce two related concepts: relative entropy and mutual information.

The relative entropy is a measure of the distance between two distributions. In statistics, it arises as an expected logarithm of the likelihood ratio. The relative entropy $D(p\|q)$ is a measure of the inefficiency of assuming that the distribution is q when the true distribution is p . For example, if we knew the true distribution of the random variable, then we could construct a code with average description length $H(p)$. If, instead, we used the code for a distribution q , we would need $H(p) + D(p\|q)$ bits on the average to describe the random variable.

Definition: The *relative entropy* or *Kullback Leibler distance* between two probability mass functions $p(x)$ and $q(x)$ is defined as

$$D(p\|q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} \quad (2.26)$$

$$= E_p \log \frac{p(X)}{q(X)}. \quad (2.27)$$

In the above definition, we use the convention (based on continuity arguments) that $0 \log \frac{0}{q} = 0$ and $p \log \frac{p}{0} = \infty$.

We will soon show that relative entropy is always non-negative and is zero if and only if $p = q$. However, it is not a true distance between distributions since it is not symmetric and does not satisfy the triangle inequality. Nonetheless, it is often useful to think of relative entropy as a “distance” between distributions.

We now introduce mutual information, which is a measure of the amount of information that one random variable contains about another random variable. It is the reduction in the uncertainty of one random variable due to the knowledge of the other.

Definition: Consider two random variables X and Y with a joint probability mass function $p(x, y)$ and marginal probability mass functions $p(x)$ and $p(y)$. The *mutual information* $I(X; Y)$ is the relative entropy between the joint distribution and the product distribution $p(x)p(y)$, i.e.,

$$I(X; Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \quad (2.28)$$

$$= D(p(x, y) \| p(x)p(y)) \quad (2.29)$$

$$= E_{p(x, y)} \log \frac{p(X, Y)}{p(X)p(Y)}. \quad (2.30)$$

Example 2.3.1: Let $\mathcal{X} = \{0, 1\}$ and consider two distributions p and q on \mathcal{X} . Let $p(0) = 1 - r$, $p(1) = r$, and let $q(0) = 1 - s$, $q(1) = s$. Then

$$D(p \| q) = (1 - r) \log \frac{1 - r}{1 - s} + r \log \frac{r}{s} \quad (2.31)$$

and

$$D(q \| p) = (1 - s) \log \frac{1 - s}{1 - r} + s \log \frac{s}{r}. \quad (2.32)$$

If $r = s$, then $D(p \| q) = D(q \| p) = 0$. If $r = 1/2$, $s = 1/4$, then we can calculate

$$D(p \| q) = \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{3}{4}} + \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{4}} = 1 - \frac{1}{2} \log 3 = 0.2075 \text{ bits}, \quad (2.33)$$

whereas

$$D(q \| p) = \frac{3}{4} \log \frac{\frac{3}{4}}{\frac{1}{2}} + \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{2}} = \frac{3}{4} \log 3 - 1 = 0.1887 \text{ bits}. \quad (2.34)$$

Note that $D(p \| q) \neq D(q \| p)$ in general.

2.4 RELATIONSHIP BETWEEN ENTROPY AND MUTUAL INFORMATION

We can rewrite the definition of mutual information $I(X; Y)$ as

$$I(X; Y) = \sum_{x, y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \quad (2.35)$$

$$= \sum_{x, y} p(x, y) \log \frac{p(x|y)}{p(x)} \quad (2.36)$$

$$= - \sum_{x, y} p(x, y) \log p(x) + \sum_{x, y} p(x, y) \log p(x|y) \quad (2.37)$$

$$= - \sum_x p(x) \log p(x) - \left(- \sum_{x, y} p(x, y) \log p(x|y) \right) \quad (2.38)$$

$$= H(X) - H(X|Y). \quad (2.39)$$

Thus the mutual information $I(X; Y)$ is the reduction in the uncertainty of X due to the knowledge of Y .

By symmetry, it also follows that

$$I(X; Y) = H(Y) - H(Y|X). \quad (2.40)$$

Thus X says as much about Y as Y says about X .

Since $H(X, Y) = H(X) + H(Y|X)$ as shown in Section 2.2, we have

$$I(X; Y) = H(X) + H(Y) - H(X, Y). \quad (2.41)$$

Finally, we note that

$$I(X; X) = H(X) - H(X|X) = H(X). \quad (2.42)$$

Thus the mutual information of a random variable with itself is the entropy of the random variable. This is the reason that entropy is sometimes referred to as *self-information*.

Collecting these results, we have the following theorem.

Theorem 2.4.1 (*Mutual information and entropy*):

$$I(X; Y) = H(X) - H(X|Y), \quad (2.43)$$

$$I(X; Y) = H(Y) - H(Y|X), \quad (2.44)$$

$$I(X; Y) = H(X) + H(Y) - H(X, Y), \quad (2.45)$$

$$I(X; Y) = I(Y; X), \quad (2.46)$$

$$I(X; X) = H(X). \quad (2.47)$$

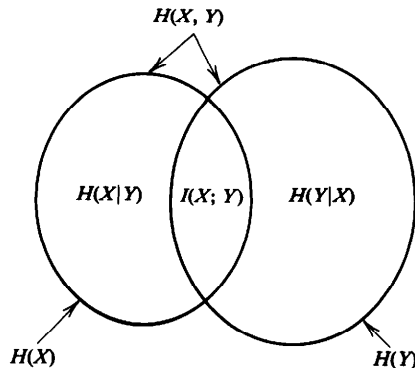


Figure 2.2. Relationship between entropy and mutual information.

The relationship between $H(X)$, $H(Y)$, $H(X, Y)$, $H(X|Y)$, $H(Y|X)$ and $I(X; Y)$ is expressed in a Venn diagram (Figure 2.2). Notice that the mutual information $I(X; Y)$ corresponds to the intersection of the information in X with the information in Y .

Example 2.4.1: For the joint distribution of example 2.2.1, it is easy to calculate the mutual information $I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) = 0.375$ bits.

2.5 CHAIN RULES FOR ENTROPY, RELATIVE ENTROPY AND MUTUAL INFORMATION

We now show that the entropy of a collection of random variables is the sum of the conditional entropies.

Theorem 2.5.1 (Chain rule for entropy): Let X_1, X_2, \dots, X_n be drawn according to $p(x_1, x_2, \dots, x_n)$. Then

$$H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1). \quad (2.48)$$

Proof: By repeated application of the two-variable expansion rule for entropies, we have

$$H(X_1, X_2) = H(X_1) + H(X_2 | X_1), \quad (2.49)$$

$$H(X_1, X_2, X_3) = H(X_1) + H(X_2, X_3 | X_1) \quad (2.50)$$

$$= H(X_1) + H(X_2 | X_1) + H(X_3 | X_2, X_1), \quad (2.51)$$

⋮

$$H(X_1, X_2, \dots, X_n) = H(X_1) + H(X_2 | X_1) + \dots + H(X_n | X_{n-1}, \dots, X_1) \quad (2.52)$$

$$= \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1). \quad (2.53)$$

Alternative Proof: We write $p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | x_{i-1}, \dots, x_1)$ and evaluate

$$\begin{aligned} & H(X_1, X_2, \dots, X_n) \\ &= - \sum_{x_1, x_2, \dots, x_n} p(x_1, x_2, \dots, x_n) \log p(x_1, x_2, \dots, x_n) \quad (2.54) \end{aligned}$$

$$= - \sum_{x_1, x_2, \dots, x_n} p(x_1, x_2, \dots, x_n) \log \prod_{i=1}^n p(x_i | x_{i-1}, \dots, x_1) \quad (2.55)$$

$$= - \sum_{x_1, x_2, \dots, x_n} \sum_{i=1}^n p(x_1, x_2, \dots, x_n) \log p(x_i | x_{i-1}, \dots, x_1) \quad (2.56)$$

$$= - \sum_{i=1}^n \sum_{x_1, x_2, \dots, x_n} p(x_1, x_2, \dots, x_n) \log p(x_i | x_{i-1}, \dots, x_1) \quad (2.57)$$

$$= - \sum_{i=1}^n \sum_{x_1, x_2, \dots, x_i} p(x_1, x_2, \dots, x_i) \log p(x_i | x_{i-1}, \dots, x_1) \quad (2.58)$$

$$= \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1). \quad \square \quad (2.59)$$

We now define the conditional mutual information as the reduction in the uncertainty of X due to knowledge of Y when Z is given.

Definition: The *conditional mutual information* of random variables X and Y given Z is defined by

$$I(X; Y|Z) = H(X|Z) - H(X|Y, Z) \quad (2.60)$$

$$= E_{p(x, y, z)} \log \frac{p(X, Y|Z)}{p(X|Z)p(Y|Z)}. \quad (2.61)$$

Mutual information also satisfies a chain rule.

Theorem 2.5.2 (*Chain rule for information*):

$$I(X_1, X_2, \dots, X_n; Y) = \sum_{i=1}^n I(X_i; Y | X_{i-1}, X_{i-2}, \dots, X_1). \quad (2.62)$$

Proof:

$$I(X_1, X_2, \dots, X_n; Y) = H(X_1, X_2, \dots, X_n) - H(X_1, X_2, \dots, X_n | Y) \quad (2.63)$$

$$\begin{aligned} &= \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1) - \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1, Y) \\ &= \sum_{i=1}^n I(X_i; Y | X_1, X_2, \dots, X_{i-1}). \quad \square \end{aligned} \quad (2.64)$$

We define a conditional version of the relative entropy.

Definition: The *conditional relative entropy* $D(p(y|x) || q(y|x))$ is the average of the relative entropies between the conditional probability mass functions $p(y|x)$ and $q(y|x)$ averaged over the probability mass function $p(x)$. More precisely,

$$D(p(y|x) || q(y|x)) = \sum_x p(x) \sum_y p(y|x) \log \frac{p(y|x)}{q(y|x)} \quad (2.65)$$

$$= E_{p(x, y)} \log \frac{p(Y|X)}{q(Y|X)}. \quad (2.66)$$

The notation for conditional relative entropy is not explicit since it omits mention of the distribution $p(x)$ of the conditioning random variable. However, it is normally understood from the context.

The relative entropy between two joint distributions on a pair of random variables can be expanded as the sum of a relative entropy and a conditional relative entropy. The chain rule for relative entropy will be used in Section 2.9 to prove a version of the second law of thermodynamics.

Theorem 2.5.3 (*Chain rule for relative entropy*):

$$D(p(x, y) \| q(x, y)) = D(p(x) \| q(x)) + D(p(y|x) \| q(y|x)). \quad (2.67)$$

Proof:

$$D(p(x, y) \| q(x, y)) = \sum_x \sum_y p(x, y) \log \frac{p(x, y)}{q(x, y)} \quad (2.68)$$

$$= \sum_x \sum_y p(x, y) \log \frac{p(x)p(y|x)}{q(x)q(y|x)} \quad (2.69)$$

$$= \sum_x \sum_y p(x, y) \log \frac{p(x)}{q(x)} + \sum_x \sum_y p(x, y) \log \frac{p(y|x)}{q(y|x)} \quad (2.70)$$

$$= D(p(x) \| q(x)) + D(p(y|x) \| q(y|x)). \quad \square \quad (2.71)$$

2.6 JENSEN'S INEQUALITY AND ITS CONSEQUENCES

In this section, we shall prove some simple properties of the quantities defined earlier. We begin with the properties of convex functions.

Definition: A function $f(x)$ is said to be *convex* over an interval (a, b) if for every $x_1, x_2 \in (a, b)$ and $0 \leq \lambda \leq 1$,

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2). \quad (2.72)$$

A function f is said to be *strictly convex* if equality holds only if $\lambda = 0$ or $\lambda = 1$.

Definition: A function f is *concave* if $-f$ is convex.

A function is convex if it always lies below any chord. A function is concave if it always lies above any chord.

Examples of convex functions include x^2 , $|x|$, e^x , $x \log x$ (for $x \geq 0$), etc. Examples of concave functions include $\log x$ and \sqrt{x} for $x \geq 0$. Figure 2.3 shows some examples of convex and concave functions. Note that linear functions $ax + b$ are both convex and concave. Convexity underlies many of the basic properties of information theoretic quantities like entropy and mutual information. Before we prove some of these properties, we derive some simple results for convex functions.

Theorem 2.6.1: *If the function f has a second derivative which is non-negative (positive) everywhere, then the function is convex (strictly convex).*

Proof: We use the Taylor series expansion of the function around x_0 , i.e.,

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x^*)}{2}(x - x_0)^2 \quad (2.73)$$

where x^* lies between x_0 and x . By hypothesis, $f''(x^*) \geq 0$, and thus the last term is always non-negative for all x .

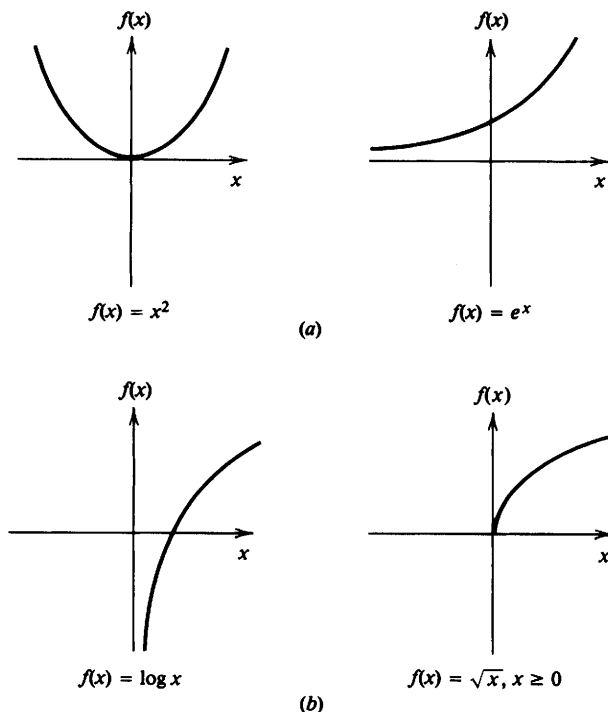


Figure 2.3. Examples of (a) convex and (b) concave functions.

We let $x_0 = \lambda x_1 + (1 - \lambda)x_2$ and take $x = x_1$ to obtain

$$f(x_1) \geq f(x_0) + f'(x_0)[(1 - \lambda)(x_1 - x_2)]. \quad (2.74)$$

Similarly, taking $x = x_2$, we obtain

$$f(x_2) \geq f(x_0) + f'(x_0)[\lambda(x_2 - x_1)]. \quad (2.75)$$

Multiplying (2.74) by λ and (2.75) by $1 - \lambda$ and adding, we obtain (2.72). The proof for strict convexity proceeds along the same lines. \square

Theorem 2.6.1 allows us to immediately verify the strict convexity of x^2 , e^x and $x \log x$ for $x \geq 0$, and the strict concavity of $\log x$ and \sqrt{x} for $x \geq 0$.

Let E denote expectation. Thus $EX = \sum_{x \in \mathcal{X}} p(x)x$ in the discrete case and $EX = \int xf(x) dx$ in the continuous case.

The next inequality is one of the most widely used in mathematics and one that underlies many of the basic results in information theory.

Theorem 2.6.2 (Jensen's inequality): *If f is a convex function and X is a random variable, then*

$$Ef(X) \geq f(EX). \quad (2.76)$$

Moreover, if f is strictly convex, then equality in (2.76) implies that $X = EX$ with probability 1, i.e., X is a constant.

Proof: We prove this for discrete distributions by induction on the number of mass points. The proof of conditions for equality when f is strictly convex will be left to the reader.

For a two mass point distribution, the inequality becomes

$$p_1 f(x_1) + p_2 f(x_2) \geq f(p_1 x_1 + p_2 x_2), \quad (2.77)$$

which follows directly from the definition of convex functions. Suppose the theorem is true for distributions with $k - 1$ mass points. Then writing $p'_i = p_i / (1 - p_k)$ for $i = 1, 2, \dots, k - 1$, we have

$$\sum_{i=1}^k p_i f(x_i) = p_k f(x_k) + (1 - p_k) \sum_{i=1}^{k-1} p'_i f(x_i) \quad (2.78)$$

$$\geq p_k f(x_k) + (1 - p_k) f\left(\sum_{i=1}^{k-1} p'_i x_i\right) \quad (2.79)$$

$$\geq f\left(p_k x_k + (1 - p_k) \sum_{i=1}^{k-1} p'_i x_i\right) \quad (2.80)$$

$$= f\left(\sum_{i=1}^k p_i x_i\right), \quad (2.81)$$

where the first inequality follows from the induction hypothesis and the second follows from the definition of convexity.

The proof can be extended to continuous distributions by continuity arguments. \square

We now use these results to prove some of the properties of entropy and relative entropy. The following theorem is of fundamental importance.

Theorem 2.6.3 (Information inequality): *Let $p(x)$, $q(x)$, $x \in \mathcal{X}$, be two probability mass functions. Then*

$$D(p\|q) \geq 0 \quad (2.82)$$

with equality if and only if

$$p(x) = q(x) \quad \text{for all } x. \quad (2.83)$$

Proof: Let $A = \{x : p(x) > 0\}$ be the support set of $p(x)$. Then

$$-D(p\|q) = -\sum_{x \in A} p(x) \log \frac{p(x)}{q(x)} \quad (2.84)$$

$$= \sum_{x \in A} p(x) \log \frac{q(x)}{p(x)} \quad (2.85)$$

$$\leq \log \sum_{x \in A} p(x) \frac{q(x)}{p(x)} \quad (2.86)$$

$$= \log \sum_{x \in A} q(x) \quad (2.87)$$

$$\leq \log \sum_{x \in \mathcal{X}} q(x) \quad (2.88)$$

$$= \log 1 \quad (2.89)$$

$$= 0, \quad (2.90)$$

where (2.86) follows from Jensen's inequality. Since $\log t$ is a strictly concave function of t , we have equality in (2.86) if and only if $q(x)/p(x) = 1$ everywhere, i.e., $p(x) = q(x)$. Hence we have $D(p\|q) = 0$ if and only if $p(x) = q(x)$ for all x . \square

Corollary (*Non-negativity of mutual information*): For any two random variables, X, Y ,

$$I(X; Y) \geq 0, \quad (2.91)$$

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

Proof: $I(X; Y) = D(p(x, y) \| p(x) p(y)) \geq 0$, with equality if and only if $p(x, y) = p(x) p(y)$, i.e., X and Y are independent. \square

Corollary:

$$D(p(y|x) \| q(y|x)) \geq 0, \quad (2.92)$$

with equality if and only if $p(y|x) = q(y|x)$ for all y and x with $p(x) > 0$.

Corollary:

$$I(X; Y|Z) \geq 0, \quad (2.93)$$

with equality if and only if X and Y are conditionally independent given Z .

We now show that the uniform distribution over the range \mathcal{X} is the maximum entropy distribution over this range. It follows that any random variable with this range has an entropy no greater than $\log |\mathcal{X}|$.

Theorem 2.6.4: $H(X) \leq \log |\mathcal{X}|$, where $|\mathcal{X}|$ denotes the number of elements in the range of X , with equality if and only if X has a uniform distribution over \mathcal{X} .

Proof: Let $u(x) = \frac{1}{|\mathcal{X}|}$ be the uniform probability mass function over \mathcal{X} , and let $p(x)$ be the probability mass function for X . Then

$$D(p \| u) = \sum p(x) \log \frac{p(x)}{u(x)} = \log |\mathcal{X}| - H(X). \quad (2.94)$$

Hence by the non-negativity of relative entropy,

$$0 \leq D(p \| u) = \log |\mathcal{X}| - H(X). \quad \square \quad (2.95)$$

Theorem 2.6.5 (*Conditioning reduces entropy*):

$$H(X|Y) \leq H(X) \quad (2.96)$$

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

Proof: $0 \leq I(X; Y) = H(X) - H(X|Y)$. \square

Intuitively, the theorem says that knowing another random variable Y can only reduce the uncertainty in X . Note that this is true only on the average. Specifically, $H(X|Y=y)$ may be greater than or less than or equal to $H(X)$, but on the average $H(X|Y) = \sum_y p(y)H(X|Y=y) \leq H(X)$. For example, in a court case, specific new evidence might increase uncertainty, but on the average evidence decreases uncertainty.

Example 2.6.1: Let (X, Y) have the following joint distribution:

$Y \backslash X$	1	2
1	0	$\frac{3}{4}$
2	$\frac{1}{8}$	$\frac{1}{8}$

Then $H(X) = H(\frac{1}{8}, \frac{7}{8}) = 0.544$ bits, $H(X|Y=1) = 0$ bits and $H(X|Y=2) = 1$ bit. We calculate $H(X|Y) = \frac{3}{4}H(X|Y=1) + \frac{1}{4}H(X|Y=2) = 0.25$ bits. Thus the uncertainty in X is increased if $Y=2$ is observed and decreased if $Y=1$ is observed, but uncertainty decreases on the average.

Theorem 2.6.6 (Independence bound on entropy): Let X_1, X_2, \dots, X_n be drawn according to $p(x_1, x_2, \dots, x_n)$. Then

$$H(X_1, X_2, \dots, X_n) \leq \sum_{i=1}^n H(X_i) \quad (2.97)$$

with equality if and only if the X_i are independent.

Proof: By the chain rule for entropies,

$$H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1) \quad (2.98)$$

$$\leq \sum_{i=1}^n H(X_i), \quad (2.99)$$

where the inequality follows directly from the previous theorem. We have equality if and only if X_i is independent of X_{i-1}, \dots, X_1 for all i , i.e., if and only if the X_i 's are independent. \square

2.7 THE LOG SUM INEQUALITY AND ITS APPLICATIONS

We now prove a simple consequence of the concavity of the logarithm, which will be used to prove some concavity results for the entropy.

Theorem 2.7.1 (*Log sum inequality*): For non-negative numbers, a_1, a_2, \dots, a_n and b_1, b_2, \dots, b_n ,

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq \left(\sum_{i=1}^n a_i \right) \log \frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i} \quad (2.100)$$

with equality if and only if $\frac{a_i}{b_i} = \text{const}$.

We again use the convention that $0 \log 0 = 0$, $a \log \frac{a}{0} = \infty$ if $a > 0$ and $0 \log \frac{0}{0} = 0$. These follow easily from continuity.

Proof: Assume without loss of generality that $a_i > 0$ and $b_i > 0$.

The function $f(t) = t \log t$ is strictly convex, since $f''(t) = \frac{1}{t} \log e > 0$ for all positive t . Hence by Jensen's inequality, we have

$$\sum \alpha_i f(t_i) \geq f\left(\sum \alpha_i t_i\right) \quad (2.101)$$

for $\alpha_i \geq 0$, $\sum_i \alpha_i = 1$. Setting $\alpha_i = b_i / \sum_{j=1}^n b_j$ and $t_i = a_i / b_i$, we obtain

$$\sum \frac{a_i}{\sum b_j} \log \frac{a_i}{b_i} \geq \sum \frac{a_i}{\sum b_j} \log \sum \frac{a_i}{\sum b_j}, \quad (2.102)$$

which is the log sum inequality. \square

We now use the log sum inequality to prove various convexity results. We begin by reproving Theorem 2.6.3, which states that $D(p\|q) \geq 0$ with equality if and only if $p(x) = q(x)$.

By the log sum inequality,

$$D(p\|q) = \sum p(x) \log \frac{p(x)}{q(x)} \quad (2.103)$$

$$\geq \left(\sum p(x) \right) \log \frac{\sum p(x)}{\sum q(x)} \quad (2.104)$$

$$= 1 \log \frac{1}{1} = 0 \quad (2.105)$$

with equality if and only if $p(x)/q(x) = c$. Since both p and q are probability mass functions, $c = 1$, and hence we have $D(p\|q) = 0$ if and only if $p(x) = q(x)$ for all x .

Theorem 2.7.2: $D(p\|q)$ is convex in the pair (p, q) , i.e., if (p_1, q_1) and (p_2, q_2) are two pairs of probability mass functions, then

$$D(\lambda p_1 + (1 - \lambda)p_2\|\lambda q_1 + (1 - \lambda)q_2) \leq \lambda D(p_1\|q_1) + (1 - \lambda)D(p_2\|q_2) \quad (2.106)$$

for all $0 \leq \lambda \leq 1$.

Proof: We apply the log sum inequality to a term on the left hand side of (2.106), i.e.,

$$\begin{aligned} (\lambda p_1(x) + (1 - \lambda)p_2(x)) \log \frac{\lambda p_1(x) + (1 - \lambda)p_2(x)}{\lambda q_1(x) + (1 - \lambda)q_2(x)} \\ \leq \lambda p_1(x) \log \frac{\lambda p_1(x)}{\lambda q_1(x)} + (1 - \lambda)p_2(x) \log \frac{(1 - \lambda)p_2(x)}{(1 - \lambda)q_2(x)}. \end{aligned} \quad (2.107)$$

Summing this over all x , we obtain the desired property. \square

Theorem 2.7.3 (Concavity of entropy): $H(p)$ is a concave function of p .

Proof:

$$H(p) = \log |\mathcal{X}| - D(p\|u), \quad (2.108)$$

where u is the uniform distribution on $|\mathcal{X}|$ outcomes. The concavity of H then follows directly from the convexity of D .

Alternative Proof: Let X_1 be a random variable with distribution p_1 taking on values in a set A . Let X_2 be another random variable with distribution p_2 on the same set. Let

$$\theta = \begin{cases} 1 & \text{with probability } \lambda \\ 2 & \text{with probability } 1 - \lambda \end{cases} \quad (2.109)$$

Let $Z = X_\theta$. Then the distribution of Z is $\lambda p_1 + (1 - \lambda)p_2$. Now since conditioning reduces entropy, we have

$$H(Z) \geq H(Z|\theta), \quad (2.110)$$

or equivalently,

$$H(\lambda p_1 + (1 - \lambda)p_2) \geq \lambda H(p_1) + (1 - \lambda)H(p_2), \quad (2.111)$$

which proves the concavity of the entropy as a function of the distribution. \square

One of the consequences of the concavity of entropy is that mixing two gases of equal entropy results in a gas with higher entropy.

Theorem 2.7.4: *Let $(X, Y) \sim p(x, y) = p(x)p(y|x)$. The mutual information $I(X; Y)$ is a concave function of $p(x)$ for fixed $p(y|x)$ and a convex function of $p(y|x)$ for fixed $p(x)$.*

Proof: To prove the first part, we expand the mutual information

$$I(X; Y) = H(Y) - H(Y|X) = H(Y) - \sum_x p(x)H(Y|X = x). \quad (2.112)$$

If $p(y|x)$ is fixed, then $p(y)$ is a linear function of $p(x)$. Hence $H(Y)$, which is a concave function of $p(y)$, is a concave function of $p(x)$. The second term is a linear function of $p(x)$. Hence the difference is a concave function of $p(x)$.

To prove the second part, we fix $p(x)$ and consider two different conditional distributions $p_1(y|x)$ and $p_2(y|x)$. The corresponding joint distributions are $p_1(x, y) = p(x) p_1(y|x)$ and $p_2(x, y) = p(x) p_2(y|x)$, and their respective marginals are $p(x)$, $p_1(y)$ and $p(x)$, $p_2(y)$. Consider a conditional distribution

$$p_\lambda(y|x) = \lambda p_1(y|x) + (1 - \lambda)p_2(y|x) \quad (2.113)$$

that is a mixture of $p_1(y|x)$ and $p_2(y|x)$. The corresponding joint distribution is also a mixture of the corresponding joint distributions,

$$p_\lambda(x, y) = \lambda p_1(x, y) + (1 - \lambda)p_2(x, y), \quad (2.114)$$

and the distribution of Y is also a mixture

$$p_\lambda(y) = \lambda p_1(y) + (1 - \lambda)p_2(y). \quad (2.115)$$

Hence if we let $q_\lambda(x, y) = p(x)p_\lambda(y)$ be the product of the marginal distributions, we have

$$q_\lambda(x, y) = \lambda q_1(x, y) + (1 - \lambda)q_2(x, y). \quad (2.116)$$

Since the mutual information is the relative entropy between the joint distribution and the product of the marginals, i.e.,

$$I(X; Y) = D(p_\lambda \| q_\lambda), \quad (2.117)$$

and relative entropy $D(p \| q)$ is a convex function of (p, q) , it follows that the mutual information is a convex function of the conditional distribution. \square

2.8 DATA PROCESSING INEQUALITY

The data processing inequality can be used to show that no clever manipulation of the data can improve the inferences that can be made from the data.

Definition: Random variables X, Y, Z are said to *form a Markov chain in that order* (denoted by $X \rightarrow Y \rightarrow Z$) if the conditional distribution of Z depends only on Y and is conditionally independent of X . Specifically, X, Y and Z form a Markov chain $X \rightarrow Y \rightarrow Z$ if the joint probability mass function can be written as

$$p(x, y, z) = p(x)p(y|x)p(z|y). \quad (2.118)$$

Some simple consequences are as follows:

- $X \rightarrow Y \rightarrow Z$ if and only if X and Z are conditionally independent given Y . Markovity implies conditional independence because

$$p(x, z|y) = \frac{p(x, y, z)}{p(y)} = \frac{p(x, y)p(z|y)}{p(y)} = p(x|y)p(z|y). \quad (2.119)$$

This is the characterization of Markov chains that can be extended to define Markov fields, which are n -dimensional random processes in which the interior and exterior are independent given the values on the boundary.

- $X \rightarrow Y \rightarrow Z$ implies that $Z \rightarrow Y \rightarrow X$. Thus the condition is sometimes written $X \leftrightarrow Y \leftrightarrow Z$.
- If $Z = f(Y)$, then $X \rightarrow Y \rightarrow Z$.

We can now prove an important and useful theorem demonstrating that no processing of Y , deterministic or random, can increase the information that Y contains about X .

Theorem 2.8.1 (Data processing inequality): If $X \rightarrow Y \rightarrow Z$, then $I(X; Y) \geq I(X; Z)$.

Proof: By the chain rule, we can expand mutual information in two different ways.

$$I(X; Y, Z) = I(X; Z) + I(X; Y|Z) \quad (2.120)$$

$$= I(X; Z) + I(X; Y|Z). \quad (2.121)$$

Since X and Z are conditionally independent given Y , we have $I(X; Z|Y) = 0$. Since $I(X; Y|Z) \geq 0$, we have

$$I(X; Y) \geq I(X; Z). \quad (2.122)$$

We have equality if and only if $I(X; Y|Z) = 0$, i.e., $X \rightarrow Z \rightarrow Y$ forms a Markov chain. Similarly, one can prove that $I(Y; Z) \geq I(X; Z)$. \square

Corollary: *In particular, if $Z = g(Y)$, we have $I(X; Y) \geq I(X; g(Y))$.*

Proof: $X \rightarrow Y \rightarrow g(Y)$ forms a Markov chain. \square

Thus functions of the data Y cannot increase the information about X .

Corollary: *If $X \rightarrow Y \rightarrow Z$, then $I(X; Y|Z) \leq I(X; Y)$.*

Proof: From (2.120) and (2.121), and using the fact that $I(X; Z|Y) = 0$ by Markovity and $I(X; Z) \geq 0$, we have

$$I(X; Y|Z) \leq I(X; Y). \quad \square \quad (2.123)$$

Thus the dependence of X and Y is decreased (or remains unchanged) by the observation of a “downstream” random variable Z .

Note that it is also possible that $I(X; Y|Z) > I(X; Y)$ when X , Y and Z do not form a Markov chain. For example, let X and Y be independent fair binary random variables, and let $Z = X + Y$. Then $I(X; Y) = 0$, but $I(X; Y|Z) = H(X|Z) - H(X|Y, Z) = H(X|Z) = P(Z = 1) H(X|Z = 1) = \frac{1}{2}$ bit.

2.9 THE SECOND LAW OF THERMODYNAMICS

One of the basic laws of physics, the second law of thermodynamics, states that the entropy of an isolated system is non-decreasing. We now explore the relationship between the second law and the entropy function that we have defined earlier in this chapter.

In statistical thermodynamics, entropy is often defined as the log of the number of microstates in the system. This corresponds exactly to our notion of entropy if all the states are equally likely. But why does the entropy increase?

We model the isolated system as a Markov chain (see Chapter 4) with transitions obeying the physical laws governing the system. Implicit in this assumption is the notion of an overall state of the system and the fact that, knowing the present state, the future of the system is independent of the past. In such a system, we can find four different interpretations of the second law. It may come as a shock to find that

the entropy does not always increase. However, *relative* entropy always decreases.

1. *Relative entropy* $D(\mu_n \parallel \mu'_n)$ decreases with n . Let μ_n and μ'_n be two probability distributions on the state space of a Markov chain at time n , and let μ_{n+1} and μ'_{n+1} be the corresponding distributions at time $n + 1$. Let the corresponding joint mass functions be denoted by p and q . Thus $p(x_n, x_{n+1}) = p(x_n) r(x_{n+1}|x_n)$ and $q(x_n, x_{n+1}) = q(x_n) r(x_{n+1}|x_n)$, where $r(\cdot | \cdot)$ is the probability transition function for the Markov chain. Then by the chain rule for relative entropy, we have two expansions:

$$\begin{aligned} D(p(x_n, x_{n+1}) \parallel q(x_n, x_{n+1})) \\ &= D(p(x_n) \parallel q(x_n)) + D(p(x_{n+1}|x_n) \parallel q(x_{n+1}|x_n)) \\ &= D(p(x_{n+1}) \parallel q(x_{n+1})) + D(p(x_n|x_{n+1}) \parallel q(x_n|x_{n+1})). \end{aligned}$$

Since both p and q are derived from the Markov chain, the conditional probability mass functions $p(x_{n+1}|x_n)$ and $q(x_{n+1}|x_n)$ are equal to $r(x_{n+1}|x_n)$ and hence $D(p(x_{n+1}|x_n) \parallel q(x_{n+1}|x_n)) = 0$. Now using the non-negativity of $D(p(x_n|x_{n+1}) \parallel q(x_n|x_{n+1}))$ (Corollary to Theorem 2.6.3), we have

$$D(p(x_n) \parallel q(x_n)) \geq D(p(x_{n+1}) \parallel q(x_{n+1})) \quad (2.124)$$

or

$$D(\mu_n \parallel \mu'_n) \geq D(\mu_{n+1} \parallel \mu'_{n+1}). \quad (2.125)$$

Consequently, the distance between the probability mass functions is decreasing with time n for any Markov chain.

An example of one interpretation of the preceding inequality is to suppose that the tax system for the redistribution of wealth is the same in Canada and in England. Then if μ_n and μ'_n represent the distributions of wealth among individuals in the two countries, this inequality shows that the relative entropy distance between the two distributions decreases with time. The wealth distributions in Canada and England will become more similar.

2. *Relative entropy* $D(\mu_n \parallel \mu)$ between a distribution μ_n on the states at time n and a stationary distribution μ decreases with n . In (2.125), μ'_n is any distribution on the states at time n . If we let μ'_n be any stationary distribution μ , then μ'_{n+1} is the same stationary distribution. Hence

$$D(\mu_n \parallel \mu) \geq D(\mu_{n+1} \parallel \mu), \quad (2.126)$$

which implies that any state distribution gets closer and closer to each stationary distribution as time passes. The sequence $D(\mu_n \parallel \mu)$ is a monotonically non-increasing non-negative sequence and must therefore have a limit. The limit is actually 0 if the stationary distribution is unique, but this is more difficult to prove.

3. *Entropy increases if the stationary distribution is uniform.* In general, the fact that the relative entropy decreases does not imply that the entropy increases. A simple counterexample is provided by any Markov chain with a non-uniform stationary distribution. If we start this Markov chain from the uniform distribution, which already is the maximum entropy distribution, the distribution will tend to the stationary distribution, which has a lower entropy than the uniform. Hence the entropy decreases with time rather than increases.

If, however, the stationary distribution is the uniform distribution, then we can express the relative entropy as

$$D(\mu_n \parallel \mu) = \log |\mathcal{X}| - H(\mu_n) = \log |\mathcal{X}| - H(X_n). \quad (2.127)$$

In this case the monotonic decrease in relative entropy implies a monotonic increase in entropy. This is the explanation that ties in most closely with statistical thermodynamics, where all the microstates are equally likely. We now characterize processes having a uniform stationary distribution.

Definition: A probability transition matrix $[P_{ij}]$, $P_{ij} = \Pr\{X_{n+1} = j | X_n = i\}$ is called *doubly stochastic* if

$$\sum_i P_{ij} = 1, \quad j = 1, 2, \dots \quad (2.128)$$

and

$$\sum_j P_{ij} = 1, \quad i = 1, 2, \dots \quad (2.129)$$

Remark: The uniform distribution is a stationary distribution of P if and only if the probability transition matrix is doubly stochastic. See Problem 1 in Chapter 4.

4. *The conditional entropy $H(X_n | X_1)$ increases with n for a stationary Markov process.* If the Markov process is stationary, then $H(X_n)$ is constant. So the entropy is nonincreasing. However, we will prove that $H(X_n | X_1)$ increases with n . Thus the conditional uncertainty of the future increases. We give two alternative proofs of this result. First, we use the properties of entropy,

$$H(X_n|X_1) \geq H(X_n|X_1, X_2) \quad (\text{conditioning reduces entropy}) \quad (2.130)$$

$$= H(X_n|X_2) \quad (\text{by Markovity}) \quad (2.131)$$

$$= H(X_{n-1}|X_1) \quad (\text{by stationarity}). \quad (2.132)$$

Alternatively, by an application of the data processing inequality to the Markov chain $X_1 \rightarrow X_{n-1} \rightarrow X_n$, we have

$$I(X_1; X_{n-1}) \geq I(X_1; X_n). \quad (2.133)$$

Expanding the mutual informations in terms of entropies, we have

$$H(X_{n-1}) - H(X_{n-1}|X_1) \geq H(X_n) - H(X_n|X_1). \quad (2.134)$$

By stationarity, $H(X_{n-1}) = H(X_n)$, and hence we have

$$H(X_{n-1}|X_1) \leq H(X_n|X_1). \quad (2.135)$$

(These techniques can also be used to show that $H(X_0|X_n)$ is increasing in n for any Markov chain. See problem 35.)

5. *Shuffles increase entropy.* If T is a shuffle (permutation) of a deck of cards and X is the initial (random) position of the cards in the deck and if the choice of the shuffle T is independent of X , then

$$H(TX) \geq H(X), \quad (2.136)$$

where TX is the permutation of the deck induced by the shuffle T on the initial permutation X . Problem 31 outlines a proof.

2.10 SUFFICIENT STATISTICS

This section is a sidelight showing the power of the data processing inequality in clarifying an important idea in statistics. Suppose we have a family of probability mass functions $\{f_\theta(x)\}$ indexed by θ , and let X be a sample from a distribution in this family. Let $T(X)$ be any statistic (function of the sample) like the sample mean or sample variance. Then $\theta \rightarrow X \rightarrow T(X)$, and by the data processing inequality, we have

$$I(\theta; T(X)) \leq I(\theta; X) \quad (2.137)$$

for any distribution on θ . However, if equality holds, no information is lost.

A statistic $T(X)$ is called sufficient for θ if it contains all the information in X about θ .

Definition: A function $T(X)$ is said to be a *sufficient statistic* relative to the family $\{f_\theta(x)\}$ if X is independent of θ given $T(X)$, i.e., $\theta \rightarrow T(X) \rightarrow X$ forms a Markov chain.

This is the same as the condition for equality in the data processing inequality,

$$I(\theta; X) = I(\theta; T(X)) \quad (2.138)$$

for all distributions on θ . Hence sufficient statistics preserve mutual information and conversely.

Here are some examples of sufficient statistics:

1. Let $X_1, X_2, \dots, X_n, X_i \in \{0, 1\}$, be an independent and identically distributed (i.i.d.) sequence of coin tosses of a coin with unknown parameter $\theta = \Pr(X_i = 1)$. Given n , the number of 1's is a sufficient statistic for θ . Here $T(X_1, X_2, \dots, X_n) = \sum_{i=1}^n X_i$. In fact, we can show that given T , all sequences having that many 1's are equally likely and independent of the parameter θ . Specifically,

$$\begin{aligned} & \Pr\left\{(X_1, X_2, \dots, X_n) = (x_1, x_2, \dots, x_n) \mid \sum_{i=1}^n X_i = k\right\} \\ &= \begin{cases} \frac{1}{\binom{n}{k}} & \text{if } \sum x_i = k, \\ 0 & \text{otherwise.} \end{cases} \end{aligned} \quad (2.139)$$

Thus $\theta \rightarrow \sum X_i \rightarrow (X_1, X_2, \dots, X_n)$ forms a Markov chain, and T is a sufficient statistic for θ .

The next two examples involve probability densities instead of probability mass functions, but the theory still applies. We define entropy and mutual information for continuous random variables in Chapter 9.

2. If X is normally distributed with mean θ and variance 1, i.e., if

$$f_\theta(x) = \frac{1}{\sqrt{2\pi}} e^{-(x-\theta)^2/2} = \mathcal{N}(\theta, 1), \quad (2.140)$$

and X_1, X_2, \dots, X_n are drawn independently according to this distribution, then a sufficient statistic for θ is $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. It can be verified that the conditional distribution of X_1, X_2, \dots, X_n , conditioned on \bar{X}_n and n does not depend on θ .

3. If $f_\theta = \text{Uniform}(\theta, \theta + 1)$, then a sufficient statistic for θ is

$$T(X_1, X_2, \dots, X_n) = (\max \{X_1, X_2, \dots, X_n\}, \min \{X_1, X_2, \dots, X_n\}). \quad (2.141)$$

The proof of this is slightly more complicated, but again one can show that the distribution of the data is independent of the parameter given the statistic T .

The minimal sufficient statistic is a sufficient statistic that is a function of all other sufficient statistics.

Definition: A statistic $T(X)$ is a *minimal sufficient statistic* relative to $\{f_\theta(x)\}$ if it is a function of every other sufficient statistic U . Interpreting this in terms of the data processing inequality, this implies that

$$\theta \rightarrow T(X) \rightarrow U(X) \rightarrow X. \quad (2.142)$$

Hence a minimal sufficient statistic maximally compresses the information about θ in the sample. Other sufficient statistics may contain additional irrelevant information. For example, for a normal distribution with mean θ , the pair of functions giving the mean of all odd samples and the mean of all even samples is a sufficient statistic, but not a minimal sufficient statistic. In the preceding examples, the sufficient statistics are also minimal.

2.11 FANO'S INEQUALITY

Suppose we know a random variable Y and we wish to guess the value of a correlated random variable X . Fano's inequality relates the probability of error in guessing the random variable X to its conditional entropy $H(X|Y)$. It will be crucial in proving the converse to Shannon's second theorem in Chapter 8. From the problems at the end of the chapter, we know that the conditional entropy of a random variable X given another random variable Y is zero if and only if X is a function of Y . Hence we can estimate X from Y with zero probability of error if and only if $H(X|Y) = 0$.

Extending this argument, we expect to be able to estimate X with a low probability of error only if the conditional entropy $H(X|Y)$ is small. Fano's inequality quantifies this idea.

Suppose we wish to estimate a random variable X with a distribution $p(x)$. We observe a random variable Y which is related to X by the conditional distribution $p(y|x)$. From Y , we calculate a function $g(Y) =$

\hat{X} , which is an estimate of X . We wish to bound the probability that $\hat{X} \neq X$. We observe that $X \rightarrow Y \rightarrow \hat{X}$ forms a Markov chain. Define the probability of error

$$P_e = \Pr\{\hat{X} \neq X\}. \quad (2.143)$$

Theorem 2.11.1 (*Fano's inequality*):

$$H(P_e) + P_e \log(|\mathcal{X}| - 1) \geq H(X|Y). \quad (2.144)$$

This inequality can be weakened to

$$1 + P_e \log|\mathcal{X}| \geq H(X|Y) \quad (2.145)$$

or

$$P_e \geq \frac{H(X|Y) - 1}{\log|\mathcal{X}|}. \quad (2.146)$$

Remark: Note that $P_e = 0$ implies that $H(X|Y) = 0$, as intuition suggests.

Proof: Define an error random variable,

$$E = \begin{cases} 1 & \text{if } \hat{X} \neq X, \\ 0 & \text{if } \hat{X} = X. \end{cases} \quad (2.147)$$

Then, using the chain rule for entropies to expand $H(E, X|Y)$ in two different ways, we have

$$H(E, X|Y) = H(X|Y) + \underbrace{H(E|X, Y)}_{=0} \quad (2.148)$$

$$= \underbrace{H(E|Y)}_{\leq H(P_e)} + \underbrace{H(X|E, Y)}_{\leq P_e \log(|\mathcal{X}| - 1)}. \quad (2.149)$$

Since conditioning reduces entropy, $H(E|Y) \leq H(E) = H(P_e)$. Now since E is a function of X and $g(Y)$, the conditional entropy $H(E|X, Y)$ is equal to 0. Also, since E is a binary-valued random variable, $H(E) = H(P_e)$. The remaining term, $H(X|E, Y)$, can be bounded as follows:

$$H(X|E, Y) = \Pr(E = 0)H(X|Y, E = 0) + \Pr(E = 1)H(X|Y, E = 1) \quad (2.150)$$

$$\leq (1 - P_e)0 + P_e \log(|\mathcal{X}| - 1), \quad (2.151)$$

since given $E = 0$, $X = g(Y)$, and given $E = 1$, we can upper bound the conditional entropy by the log of the number of remaining outcomes ($|\mathcal{X}| - 1$ if $g(Y) \in \mathcal{X}$, else $|\mathcal{X}|$). Combining these results, we obtain Fano's inequality. \square

Remark: Suppose that there is no knowledge of Y . Thus X must be guessed without any information. Let $X \in \{1, 2, \dots, m\}$ and $p_1 \geq p_2 \geq \dots \geq p_m$. Then the best guess of X is $\hat{X} = 1$ and the resulting probability of error is $P_e = 1 - p_1$. Fano's inequality becomes

$$H(P_e) + P_e \log(m - 1) \geq H(X). \quad (2.152)$$

The probability mass function

$$(p_1, p_2, \dots, p_m) = \left(1 - P_e, \frac{P_e}{m - 1}, \dots, \frac{P_e}{m - 1}\right) \quad (2.153)$$

achieves this bound with equality. Thus Fano's inequality is sharp.

The following telegraphic summary omits qualifying conditions.

SUMMARY OF CHAPTER 2

Definition: The *entropy* $H(X)$ of a discrete random variable X is defined by

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log p(x). \quad (2.154)$$

Properties of H :

1. $H(X) \geq 0$.
2. $H_b(X) = (\log_b a) H_a(X)$.
3. (*Conditioning reduces entropy*) For any two random variables, X and Y , we have

$$H(X|Y) \leq H(X) \quad (2.155)$$

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

4. $H(X_1, X_2, \dots, X_n) \leq \sum_{i=1}^n H(X_i)$, with equality if and only if the random variables X_i are independent.
5. $H(X) \leq \log |\mathcal{X}|$ with equality if and only if X is uniformly distributed over \mathcal{X} .
6. $H(p)$ is concave in p .

Definition: The *relative entropy* $D(p||q)$ of the probability mass function p with respect to the probability mass function q is defined by

$$D(p\|q) = \sum_x p(x) \log \frac{p(x)}{q(x)}. \quad (2.156)$$

Definition: The *mutual information* between two random variables X and Y is defined as

$$I(X; Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}. \quad (2.157)$$

Alternative expressions:

$$H(X) = E_p \log \frac{1}{p(X)} \quad (2.158)$$

$$H(X, Y) = E_p \log \frac{1}{p(X, Y)} \quad (2.159)$$

$$H(X|Y) = E_p \log \frac{1}{p(X|Y)} \quad (2.160)$$

$$I(X; Y) = E_p \log \frac{p(X, Y)}{p(X)p(Y)} \quad (2.161)$$

$$D(p\|q) = E_p \log \frac{p(X)}{q(X)} \quad (2.162)$$

Properties of D and I :

1. $I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) = H(X) + H(Y) - H(X, Y)$.
2. $D(p\|q) \geq 0$ with equality if and only if $p(x) = q(x)$, for all $x \in \mathcal{X}$.
3. $I(X; Y) = D(p(x, y)\|p(x)p(y)) \geq 0$, with equality if and only if $p(x, y) = p(x)p(y)$, i.e., X and Y are independent.
4. If $|\mathcal{X}| = m$, and u is the uniform distribution over \mathcal{X} , then $D(p\|u) = \log m - H(p)$.
5. $D(p\|q)$ is convex in the pair (p, q) .

Chain rules

Entropy: $H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i|X_{i-1}, \dots, X_1)$.

Mutual

information: $I(X_1, X_2, \dots, X_n; Y) = \sum_{i=1}^n I(X_i; Y|X_1, X_2, \dots, X_{i-1})$.

Relative entropy: $D(p(x, y)\|q(x, y)) = D(p(x)\|q(x)) + D(p(y|x)\|q(y|x))$.

Jensen's inequality: If f is a convex function, then $Ef(X) \geq f(EX)$.

Log sum inequality: For n positive numbers, a_1, a_2, \dots, a_n and b_1, b_2, \dots, b_n ,

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq \left(\sum_{i=1}^n a_i \right) \log \frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i} \quad (2.163)$$

with equality if and only if $a_i/b_i = \text{constant}$.

Data processing inequality: If $X \rightarrow Y \rightarrow Z$ forms a Markov chain, then $I(X; Y) \geq I(X; Z)$

Second law of thermodynamics: For a Markov chain,

1. Relative entropy $D(\mu_n \| \mu'_n)$ decreases with time.
2. Relative entropy $D(\mu_n \| \mu)$ between a distribution and the stationary distribution decreases with time.
3. Entropy $H(X_n)$ increases if the stationary distribution is uniform.
4. The conditional entropy $H(X_n | X_1)$ increases with time for a stationary Markov chain.
5. The conditional entropy $H(X_0 | X_n)$ of the initial condition X_0 increases for any Markov chain.

Sufficient statistic: $T(X)$ is sufficient relative to $\{f_\theta(x)\}$ if and only if $I(\theta; X) = I(\theta; T(X))$ for all distributions on θ .

Fano's inequality: Let $P_e = \Pr\{g(Y) \neq X\}$, where g is any function of Y . Then

$$H(P_e) + P_e \log(|\mathcal{X}| - 1) \geq H(X|Y). \quad (2.164)$$

PROBLEMS FOR CHAPTER 2

1. *Coin flips.* A fair coin is flipped until the first head occurs. Let X denote the number of flips required.
 - (a) Find the entropy $H(X)$ in bits. The following expressions may be useful:

$$\sum_{n=1}^{\infty} r^n = \frac{r}{1-r}, \quad \sum_{n=1}^{\infty} nr^n = \frac{r}{(1-r)^2}.$$

- (b) A random variable X is drawn according to this distribution. Find an "efficient" sequence of yes-no questions of the form, "Is X contained in the set S ?" Compare $H(X)$ to the expected number of questions required to determine X .
2. *Entropy of functions.* 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
 - (a) $Y = 2^X$?
 - (b) $Y = \cos X$?
3. *Minimum entropy.* What is the minimum value of $H(p_1, \dots, p_n) = H(\mathbf{p})$ as \mathbf{p} ranges over the set of n -dimensional probability vectors? Find all \mathbf{p} 's which achieve this minimum.
4. *Axiomatic definition of entropy.* If we assume certain axioms for our measure of information, then we will be forced to use a logarithmic

measure like entropy. Shannon used this to justify his initial definition of entropy. In this book, we will rely more on the other properties of entropy rather than its axiomatic derivation to justify its use. The following problem is considerably more difficult than the other problems in this section.

If a sequence of symmetric functions $H_m(p_1, p_2, \dots, p_m)$ satisfies the following properties,

- Normalization: $H_2(\frac{1}{2}, \frac{1}{2}) = 1$,
- Continuity: $H_2(p, 1-p)$ is a continuous function of p ,
- Grouping: $H_m(p_1, p_2, \dots, p_m) = H_{m-1}(p_1 + p_2, p_3, \dots, p_m) + (p_1 + p_2) H_2(\frac{p_1}{p_1+p_2}, \frac{p_2}{p_1+p_2})$,

prove that H_m must be of the form

$$H_m(p_1, p_2, \dots, p_m) = - \sum_{i=1}^m p_i \log p_i, \quad m = 2, 3, \dots \tag{2.165}$$

There are various other axiomatic formulations which also result in the same definition of entropy. See, for example, the book by Csiszar and Körner [83].

5. *Entropy of functions of a random variable.* Let X be a discrete random variable. Show that the entropy of a function of X is less than or equal to the entropy of X by justifying the following steps:

$$H(X, g(X)) \stackrel{(a)}{=} H(X) + H(g(X)|X) \tag{2.166}$$

$$\stackrel{(b)}{=} H(X); \tag{2.167}$$

$$H(X, g(X)) \stackrel{(c)}{=} H(g(X)) + H(X|g(X)) \tag{2.168}$$

$$\stackrel{(d)}{\geq} H(g(X)). \tag{2.169}$$

Thus $H(g(X)) \leq H(X)$.

6. *Zero conditional entropy.* Show that if $H(Y|X) = 0$, then Y is a function of X , i.e., for all x with $p(x) > 0$, there is only one possible value of y with $p(x, y) > 0$.
7. *Pure randomness and bent coins.* Let X_1, X_2, \dots, X_n denote the outcomes of independent flips of a *bent* coin. Thus $\Pr\{X_i = 1\} = p$, $\Pr\{X_i = 0\} = 1 - p$, where p is unknown. We wish to obtain a sequence Z_1, Z_2, \dots, Z_K of *fair* coin flips from X_1, X_2, \dots, X_n . Toward this end let $f: \mathcal{X}^n \rightarrow \{0, 1\}^*$ (where $\{0, 1\}^* = \{\Lambda, 0, 1, 00, 01, \dots\}$ is the set of all finite length binary sequences) be a mapping $f(X_1, X_2, \dots, X_n) = (Z_1, Z_2, \dots, Z_K)$, where $Z_i \sim \text{Bernoulli}(\frac{1}{2})$, and K may depend on (X_1, \dots, X_n) . In order that the sequence Z_1, Z_2, \dots

appear to be fair coin flips, the map f from bent coin flips to fair flips must have the property that all 2^k sequences (Z_1, Z_2, \dots, Z_k) of a given length k have equal probability (possibly 0), for $k = 1, 2, \dots$. For example, for $n = 2$, the map $f(01) = 0$, $f(10) = 1$, $f(00) = f(11) = \Lambda$ (the null string), has the property that $\Pr\{Z_1 = 1|K = 1\} = \Pr\{Z_1 = 0|K = 1\} = \frac{1}{2}$.

Give reasons for the following inequalities:

$$\begin{aligned} nH(p) &\stackrel{(a)}{=} H(X_1, \dots, X_n) \\ &\stackrel{(b)}{\geq} H(Z_1, Z_2, \dots, Z_K, K) \\ &\stackrel{(c)}{=} H(K) + H(Z_1, \dots, Z_K|K) \\ &\stackrel{(d)}{=} H(K) + E(K) \\ &\stackrel{(e)}{\geq} EK. \end{aligned}$$

Thus no more than $nH(p)$ fair coin tosses can be derived from (X_1, \dots, X_n) , on the average.

(f) Exhibit a good map f on sequences of length 4.

8. *World Series.* The World Series is a seven-game series that terminates as soon as either team wins four games. Let X be the random variable that represents the outcome of a World Series between teams A and B; possible values of X are AAAA, BABABAB, and BBBAAAA. Let Y be the number of games played, which ranges from 4 to 7. Assuming that A and B are equally matched and that the games are independent, calculate $H(X)$, $H(Y)$, $H(Y|X)$, and $H(X|Y)$.
9. *Infinite entropy.* 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 $\Pr(X = n) = (An \log^2 n)^{-1}$ for $n = 2, 3, \dots$ has $H(X) = +\infty$.
10. *Conditional mutual information vs. unconditional mutual information.* Give examples of joint random variables X , Y and Z such that
 - (a) $I(X; Y|Z) < I(X; Y)$,
 - (b) $I(X; Y|Z) > I(X; Y)$.
11. *Average entropy.* Let $H(p) = -p \log_2 p - (1-p) \log_2(1-p)$ be the binary entropy function.
 - (a) Evaluate $H(1/4)$ using the fact that $\log_2 3 \approx 1.584$. *Hint:* Consider an experiment with four equally likely outcomes, one of which is more interesting than the others.
 - (b) Calculate the average entropy $H(p)$ when the probability p is chosen uniformly in the range $0 \leq p \leq 1$.

- (c) (*Optional*) Calculate the average entropy $H(p_1, p_2, p_3)$ where (p_1, p_2, p_3) is a uniformly distributed probability vector. Generalize to dimension n .

12. *Venn diagrams.* Using Venn diagrams, we can see that the mutual information common to three random variables X, Y and Z should be defined by

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

This quantity is symmetric in X, Y and Z , despite the preceding asymmetric definition. Unfortunately, $I(X; Y; Z)$ is not necessarily nonnegative. Find X, Y and Z such that $I(X; Y; Z) < 0$, and prove the following two identities:

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

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

The first identity can be understood using the Venn diagram analogy for entropy and mutual information. The second identity follows easily from the first.

13. *Coin weighing.* Suppose one has n coins, among which there may or may not be one counterfeit coin. If there is a counterfeit coin, it may be either heavier or lighter than the other coins. The coins are to be weighed by a balance.
- (a) Find an upper bound on the number of coins n so that k weighings will find the counterfeit coin (if any) and correctly declare it to be heavier or lighter.
- (b) (*Difficult*) What is the coin weighing strategy for $k = 3$ weighings and 12 coins?
14. *Drawing with and without replacement.* An urn contains r red, w white, and b black balls. Which has higher entropy, drawing $k \geq 2$ balls from the urn with replacement or without replacement? Set it up and show why. (There is both a hard way and a relatively simple way to do this.)
15. *A metric.* A function $\rho(x, y)$ is a metric if for all x, y ,
- $\rho(x, y) \geq 0$
 - $\rho(x, y) = \rho(y, x)$
 - $\rho(x, y) = 0$ if and only if $x = y$
 - $\rho(x, y) + \rho(y, z) \geq \rho(x, z)$.

- (a) Show that $\rho(X, Y) = H(X|Y) + H(Y|X)$ satisfies the first, second and fourth properties above. If we say that $X = Y$ if there is a one-to-one function mapping X to Y , then the third property is also satisfied, and $\rho(X, Y)$ is a metric.
- (b) Verify that $\rho(X, Y)$ can also be expressed as

$$\rho(X, Y) = H(X) + H(Y) - 2I(X; Y) \quad (2.170)$$

$$= H(X, Y) - I(X; Y) \quad (2.171)$$

$$= 2H(X, Y) - H(X) - H(Y). \quad (2.172)$$

16. *Example of joint entropy.* Let $p(x, y)$ be given by

		Y	
		0	1
X	0	$\frac{1}{3}$	$\frac{1}{3}$
	1	0	$\frac{1}{3}$

Find

- (a) $H(X)$, $H(Y)$.
- (b) $H(X|Y)$, $H(Y|X)$.
- (c) $H(X, Y)$.
- (d) $H(Y) - H(Y|X)$.
- (e) $I(X; Y)$.
- (f) Draw a Venn diagram for the quantities in (a) through (e).
17. *Inequality.* Show $\ln x \geq 1 - \frac{1}{x}$ for $x > 0$.
18. *Entropy of a sum.* Let X and Y be random variables that take on values x_1, x_2, \dots, x_r and y_1, y_2, \dots, y_s , respectively. Let $Z = X + Y$.
- (a) Show that $H(Z|X) = H(Y|X)$. Argue that if X, Y are independent, then $H(Y) \leq H(Z)$ and $H(X) \leq H(Z)$. Thus the addition of *independent* random variables adds uncertainty.
- (b) Give an example (of necessarily dependent random variables) in which $H(X) > H(Z)$ and $H(Y) > H(Z)$.
- (c) Under what conditions does $H(Z) = H(X) + H(Y)$?
19. *Entropy of a disjoint mixture.* Let X_1 and X_2 be discrete random variables drawn according to probability mass functions $p_1(\cdot)$ and $p_2(\cdot)$ over the respective alphabets $\mathcal{X}_1 = \{1, 2, \dots, m\}$ and $\mathcal{X}_2 = \{m+1, \dots, n\}$. Let

$$X = \begin{cases} X_1, & \text{with probability } \alpha, \\ X_2, & \text{with probability } 1 - \alpha. \end{cases}$$

- (a) Find $H(X)$ in terms of $H(X_1)$ and $H(X_2)$ and α .
- (b) Maximize over α to show that $2^{H(X)} \leq 2^{H(X_1)} + 2^{H(X_2)}$ and interpret using the notion that $2^{H(X)}$ is the effective alphabet size.

20. *A measure of correlation.* Let X_1 and X_2 be identically distributed, but not necessarily independent. Let

$$\rho = 1 - \frac{H(X_2|X_1)}{H(X_1)}.$$

- (a) Show $\rho = I(X_1; X_2)/H(X_1)$.
 (b) Show $0 \leq \rho \leq 1$.
 (c) When is $\rho = 0$?
 (d) When is $\rho = 1$?
21. *Data processing.* Let $X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow \dots \rightarrow X_n$ form a Markov chain in this order; i.e., 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.

22. *Bottleneck.* Suppose a (non-stationary) Markov chain starts in one of n states, necks down to $k < n$ states, and then fans back to $m > k$ states. Thus $X_1 \rightarrow X_2 \rightarrow X_3$, $X_1 \in \{1, 2, \dots, n\}$, $X_2 \in \{1, 2, \dots, k\}$, $X_3 \in \{1, 2, \dots, m\}$.
- (a) Show that the dependence of X_1 and X_3 is limited by the bottleneck by proving that $I(X_1; X_3) \leq \log k$.
 (b) Evaluate $I(X_1; X_3)$ for $k = 1$, and conclude that no dependence can survive such a bottleneck.
23. *Run length coding.* Let X_1, X_2, \dots, X_n be (possibly dependent) binary random variables. Suppose one calculates the run lengths $\mathbf{R} = (R_1, R_2, \dots)$ of this sequence (in order as they occur). For example, the sequence $\mathbf{x} = 0001100100$ yields run lengths $\mathbf{R} = (3, 2, 2, 1, 2)$. Compare $H(X_1, X_2, \dots, X_n)$, $H(\mathbf{R})$ and $H(X_n, \mathbf{R})$. Show all equalities and inequalities, and bound all the differences.
24. *Markov's inequality for probabilities.* Let $p(x)$ be a probability mass function. Prove, for all $d \geq 0$,

$$\Pr\{p(X) \leq d\} \log\left(\frac{1}{d}\right) \leq H(X). \quad (2.173)$$

25. *Logical order of ideas.* Ideas have been developed in order of need, and then generalized if necessary. Reorder the following ideas, strongest first, implications following:
- (a) Chain rule for $I(X_1, \dots, X_n; Y)$, chain rule for $D(p(x_1, \dots, x_n) \| q(x_1, x_2, \dots, x_n))$, and chain rule for $H(X_1, X_2, \dots, X_n)$.
 (b) $D(f \| g) \geq 0$, Jensen's inequality, $I(X; Y) \geq 0$.
26. *Second law of thermodynamics.* Let $X_1, X_2, X_3 \dots$ be a stationary first-order Markov chain. In Section 2.9, it was shown that $H(X_n|X_1) \geq H(X_{n-1}|X_1)$ for $n = 2, 3, \dots$. Thus conditional uncertainty about the future grows with time. This is true although the unconditional

uncertainty $H(X_n)$ remains constant. However, show by example that $H(X_n|X_1 = x_1)$ does not necessarily grow with n for every x_1 .

27. *Conditional mutual information.* Consider a sequence of n binary random variables X_1, X_2, \dots, X_n . Each sequence with an even number of 1's has probability $2^{-(n-1)}$ and each sequence with an odd number of 1's has probability 0. Find the mutual informations

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

28. *Mixing increases entropy.* Show that the entropy of the probability distribution, $(p_1, \dots, p_i, \dots, p_j, \dots, p_m)$, is less than the entropy of the distribution $(p_1, \dots, \frac{p_i + p_j}{2}, \dots, \frac{p_i + p_j}{2}, \dots, p_m)$. Show that in general any transfer of probability that makes the distribution more uniform increases the entropy.
29. *Inequalities.* Let X, Y and Z be joint random variables. Prove the following inequalities and find conditions for equality.
- $H(X, Y|Z) \geq H(X|Z)$.
 - $I(X, Y; Z) \geq I(X; Z)$.
 - $H(X, Y, Z) - H(X, Y) \leq H(X, Z) - H(X)$.
 - $I(X; Z|Y) \geq I(Z; Y|X) - I(Z; Y) + I(X; Z)$.
30. *Maximum entropy.* Find the probability mass function $p(x)$ that maximizes the entropy $H(X)$ of a non-negative integer-valued random variable X subject to the constraint

$$EX = \sum_{n=0}^{\infty} np(n) = A$$

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

31. *Shuffles increase entropy.* Argue that for any distribution on shuffles T and any distribution on card positions X that

$$H(TX) \geq H(TX|T) \tag{2.174}$$

$$= H(T^{-1}TX|T) \tag{2.175}$$

$$= H(X|T) \tag{2.176}$$

$$= H(X), \tag{2.177}$$

if X and T are independent.

32. *Conditional entropy.* Under what conditions does $H(X|g(Y)) = H(X|Y)$?
33. *Fano's inequality.* Let $\Pr(X = i) = p_i, i = 1, 2, \dots, m$ and let $p_1 \geq p_2 \geq p_3 \geq \dots \geq p_m$. The minimal probability of error predictor of X is $\hat{X} = 1$, with resulting probability of error $P_e = 1 - p_1$. Maximize $H(\mathbf{p})$ subject to the constraint $1 - p_1 = P_e$ to find a bound on P_e in terms of H . This is Fano's inequality in the absence of conditioning.

34. *Monotonic convergence of the empirical distribution.* Let \hat{p}_n denote the empirical probability mass function corresponding to X_1, X_2, \dots, X_n i.i.d. $\sim p(x)$, $x \in \mathcal{X}$. Specifically,

$$\hat{p}_n(x) = \frac{1}{n} \sum_{i=1}^n I(X_i = x) \quad (2.178)$$

is the proportion of times that $X_i = x$ in the first n samples, where I is an indicator function.

- (a) Show for \mathcal{X} binary that

$$ED(\hat{p}_{2n} \| p) \leq ED(\hat{p}_n \| p). \quad (2.179)$$

Thus the expected relative entropy “distance” from the empirical distribution to the true distribution decreases with sample size.

Hint: Write $\hat{p}_{2n} = \frac{1}{2} \hat{p}_n + \frac{1}{2} \hat{p}'_n$ and use the convexity of D .

- (b) Show for an arbitrary discrete \mathcal{X} that

$$ED(\hat{p}_n \| p) \leq ED(\hat{p}_{n-1} \| p). \quad (2.180)$$

35. *Entropy of initial conditions.* Prove that $H(X_0|X_n)$ is non-decreasing with n for any Markov chain.

HISTORICAL NOTES

The concept of entropy was first introduced in thermodynamics, where it was used to provide a statement of the second law of thermodynamics. Later, statistical mechanics provided a connection between the macroscopic property of entropy and the microscopic state of the system. This work was the crowning achievement of Boltzmann, who had the equation $S = k \ln W$ inscribed as the epitaph on his gravestone.

In the 1930s, Hartley introduced a logarithmic measure of information for communication. His measure was essentially the logarithm of the alphabet size. Shannon [238] was the first to define entropy and mutual information as defined in this chapter. Relative entropy was first defined by Kullback and Leibler [167]. It is known under a variety of names, including the Kullback Leibler distance, cross entropy, information divergence and information for discrimination, and has been studied in detail by Csiszár [78] and Amari [10].

Many of the simple properties of these quantities were developed by Shannon. Fano’s inequality was proved in Fano [105]. The notion of sufficient statistic was defined by Fisher [111], and the notion of the minimal sufficient statistic was introduced by Lehmann and Scheffé [174]. The relationship of mutual information and sufficiency is due to Kullback [165].

The relationship between information theory and thermodynamics has been discussed extensively by Brillouin [46] and Jaynes [143]. Although the basic theorems of information theory were originally derived for a communication system, attempts have been made to compare these theorems with the fundamental laws of physics. There have also been attempts to determine whether there are any fundamental physical limits to computation, including work by Bennett [24] and Bennett and Landauer [25].