

2 Fundamentals of information theory

We begin with a brief overview of some of the fundamental concepts and mathematical tools of information theory. This allows us to establish notation and to set the stage for the results presented in subsequent chapters. For a comprehensive introduction to the fundamental concepts and methods of information theory, we refer the interested reader to the textbooks of Gallager [2], Cover and Thomas [3], Yeung [4], and Csiszár and Körner [5].

The rest of the chapter is organized as follows. Section 2.1 provides an overview of the basic mathematical tools and metrics that are relevant for subsequent chapters. Section 2.2 illustrates the fundamental proof techniques used in information theory by discussing the point-to-point communication problem and Shannon's coding theorems. Section 2.3 is entirely devoted to network information theory, with a special emphasis on distributed source coding and multi-user communications as they relate to information-theoretic security.

2.1 Mathematical tools of information theory

The following subsections describe a powerful set of metrics and tools that are useful to characterize the fundamental limits of communication systems. All results are stated without proof through a series of lemmas and theorems, and we refer the reader to standard textbooks [2, 3, 4] for details. Unless specified otherwise, all random variables and random vectors used throughout this book are real-valued random vectors.

2.1.1 Useful bounds

We start by recalling a few inequalities that are useful to bound the probabilities of rare events.

Lemma 2.1 (Markov's inequality). *Let X be a non-negative real-valued random variable. Then,*

$$\forall a > 0 \quad \mathbb{P}[X \geq a] \leq \frac{\mathbb{E}_X[X]}{a}.$$

The following consequence of Markov's inequality is particularly useful.

Lemma 2.2 (Selection lemma). *Let $X_n \in \mathcal{X}_n$ be a random variable and let \mathcal{F} be a finite set of functions $f : \mathcal{X}_n \rightarrow \mathbb{R}^+$ such that $|\mathcal{F}|$ does not depend on n and*

$$\forall f \in \mathcal{F} \quad \mathbb{E}_{X_n}[f(X_n)] \leq \delta(n).$$

Then, there exists a specific realization x_n of X_n such that

$$\forall f \in \mathcal{F} \quad f(x_n) \leq \delta(n).$$

Proof. Let $\epsilon_n \triangleq \delta(n)$. Using the union bound and Markov's inequality, we obtain

$$\begin{aligned} \mathbb{P}_{X_n} \left[\bigcup_{f \in \mathcal{F}} \{f(X_n) \geq (|\mathcal{F}| + 1)\epsilon_n\} \right] &\leq \sum_{f \in \mathcal{F}} \mathbb{P}_{X_n}[f(X_n) \geq (|\mathcal{F}| + 1)\epsilon_n] \\ &\leq \sum_{f \in \mathcal{F}} \frac{\mathbb{E}_{X_n}[f(X_n)]}{(|\mathcal{F}| + 1)\epsilon_n} \\ &\leq \frac{|\mathcal{F}|}{|\mathcal{F}| + 1} \\ &< 1. \end{aligned}$$

Therefore, there exists *at least* one realization x_n of X_n such that

$$\forall f \in \mathcal{F} \quad f(x_n) \leq (|\mathcal{F}| + 1)\epsilon_n.$$

Since $\epsilon_n = \delta(n)$ and $|\mathcal{F}|$ is finite and independent of n , we can write $(|\mathcal{F}| + 1)\epsilon_n$ as $\delta(n)$. \square

We call Lemma 2.2 the “selection lemma” because it tells us that, if $\mathbb{E}_{X_n}[f(X_n)] \leq \delta(n)$ for all $f \in \mathcal{F}$, we can select a specific realization x_n such that $f(x_n) \leq \delta(n)$ for all $f \in \mathcal{F}$. Two other useful consequences of Markov's inequality are Chebyshev's inequality and Chernov bounds.

Lemma 2.3 (Chebyshev's inequality). *Let X be a real-valued random variable. Then,*

$$\forall a > 0 \quad \mathbb{P}[|X - \mathbb{E}_X[X]| \geq a] \leq \frac{\text{Var}(X)}{a^2}.$$

Lemma 2.4 (Chernov bounds). *Let X be a real-valued random variable. Then, for all $a > 0$,*

$$\begin{aligned} \forall s > 0 \quad \mathbb{P}[X \geq a] &\leq \mathbb{E}_X[e^{sX}]e^{-sa}, \\ \forall s < 0 \quad \mathbb{P}[X \leq a] &\leq \mathbb{E}_X[e^{sX}]e^{-sa}. \end{aligned}$$

2.1.2 Entropy and mutual information

In this section, we define a series of useful information-theoretic quantities, whose operational significance will become clear in the next section.

Definition 2.1. Let X and X' be two discrete random variables defined on the same alphabet \mathcal{X} . The variational distance between X and X' is

$$\mathbb{V}(X, X') \triangleq \sum_{x \in \mathcal{X}} |p_X(x) - p_{X'}(x)|.$$

Definition 2.2. Let $X \in \mathcal{X}$ be a discrete random variable with distribution p_X . The Shannon entropy (or entropy for short) of X is defined as

$$\mathbb{H}(X) \triangleq - \sum_{x \in \mathcal{X}} p_X(x) \log p_X(x),$$

with the convention that $0 \log 0 \triangleq 0$. Unless specified otherwise, all logarithms are taken to the base two and the unit for entropy is called a bit.

If $\mathcal{X} = \{0, 1\}$, then X is a binary random variable and its entropy depends solely on the parameter $p = \mathbb{P}[X = 0]$. The binary entropy function is defined as

$$\mathbb{H}_b(p) \triangleq -p \log p - (1 - p) \log(1 - p).$$

Lemma 2.5. For any discrete random variable $X \in \mathcal{X}$

$$0 \leq \mathbb{H}(X) \leq \log |\mathcal{X}|.$$

The equality $\mathbb{H}(X) = 0$ holds if and only if X is a constant while the equality $\mathbb{H}(X) = \log |\mathcal{X}|$ holds if and only if X is uniform on \mathcal{X} .

Conceptually, $\mathbb{H}(X)$ can be viewed as a measure of the average amount of information contained in X or, equivalently, the amount of uncertainty that subsists until the outcome of X is revealed.

Proposition 2.1 (Csiszár and Körner). Let X and X' be two discrete random variables defined on the same alphabet \mathcal{X} . Then,

$$|\mathbb{H}(X) - \mathbb{H}(X')| \leq \mathbb{V}(X, X') \log \left(\frac{|\mathcal{X}|}{\mathbb{V}(X, X')} \right).$$

Definition 2.3. Let $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$ be two discrete random variables with joint distribution p_{XY} . The joint entropy of X and Y is defined as

$$\mathbb{H}(XY) \triangleq - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{XY}(x, y) \log p_{XY}(x, y).$$

Definition 2.4. Let $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$ be two discrete random variables with joint distribution p_{XY} . The conditional entropy of X given Y is defined as

$$\mathbb{H}(Y|X) \triangleq - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{XY}(x, y) \log p_{Y|X}(y|x).$$

By expanding $\mathbb{H}(XY)$ with Bayes' rule, one can verify that

$$\mathbb{H}(XY) = \mathbb{H}(X) + \mathbb{H}(Y|X).$$

This expansion generalizes to the entropy of a random vector $X^n = (X_1, \dots, X_n)$ as

$$\begin{aligned}\mathbb{H}(X^n) &= \mathbb{H}(X_1) + \mathbb{H}(X_2|X_1) + \dots + \mathbb{H}(X_n|X^{n-1}) \\ &= \sum_{i=1}^n \mathbb{H}(X_i|X^{i-1}),\end{aligned}$$

with the convention that $\mathbb{H}(X_1|X^0) \triangleq \mathbb{H}(X_1)$. This expansion is known as the *chain rule of entropy*.

Lemma 2.6 (“Conditioning does not increase entropy”). *Let $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$ be two discrete random variables with joint distribution p_{XY} . Then,*

$$\mathbb{H}(X|Y) \leq \mathbb{H}(X).$$

In other words, this lemma asserts that knowledge of Y cannot increase our uncertainty about X .

Definition 2.5. *Let $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$ be two discrete random variables with joint distribution p_{XY} . The mutual information between X and Y is defined as*

$$\mathbb{I}(X; Y) \triangleq \mathbb{H}(X) - \mathbb{H}(X|Y).$$

Let $X \in \mathcal{X}$, $Y \in \mathcal{Y}$, and $Z \in \mathcal{Z}$ be discrete random variables with joint distribution p_{XYZ} . The conditional mutual information between X and Y given Z is

$$\mathbb{I}(X; Y|Z) \triangleq \mathbb{H}(X|Z) - \mathbb{H}(X|YZ).$$

Intuitively, $\mathbb{I}(X; Y)$ represents the uncertainty about X that is not resolved by the observation of Y . By using the chain rule of entropy, one can expand the mutual information between a random vector $X^n = (X_1, \dots, X_n)$ and a random variable Y as

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

with the convention that $\mathbb{I}(X_1; Y|X^0) \triangleq \mathbb{I}(X_1; Y)$. This expansion is known as the *chain rule of mutual information*.

Lemma 2.7. *Let $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$ be two discrete random variables with joint distribution p_{XY} . Then,*

$$0 \leq \mathbb{I}(X; Y) \leq \min(\mathbb{H}(X), \mathbb{H}(Y)).$$

The equality $\mathbb{I}(X; Y) = 0$ holds if and only if X and Y are independent. The equality $\mathbb{I}(X; Y) = \mathbb{H}(X)$ ($\mathbb{I}(X; Y) = \mathbb{H}(Y)$) holds if and only if X is a function of Y (Y is a function of X).

Lemma 2.8. *Let $X \in \mathcal{X}$, $Y \in \mathcal{Y}$, and $Z \in \mathcal{Z}$ be three discrete random variables with joint distribution p_{XYZ} . Then,*

$$0 \leq \mathbb{I}(X; Y|Z) \leq \min(\mathbb{H}(X|Z), \mathbb{H}(Y|Z)).$$

The equality $\mathbb{I}(X; Y|Z) = 0$ holds if and only if X and Y are conditionally independent given Z . In this case, we say that $X \rightarrow Z \rightarrow Y$ forms a Markov chain. The equality $\mathbb{I}(X; Y|Z) = \mathbb{H}(X|Z)$ ($\mathbb{I}(X; Y|Z) = \mathbb{H}(Y|Z)$) holds if and only if X is a function of Y and Z (Y is a function of X and Z).

Lemma 2.9 (Data-processing inequality). *Let $X \in \mathcal{X}$, $Y \in \mathcal{Y}$, and $Z \in \mathcal{Z}$ be three discrete random variables such that $X \rightarrow Y \rightarrow Z$ forms a Markov chain. Then,*

$$\mathbb{I}(X; Y) \geq \mathbb{I}(X; Z).$$

An equivalent form of the data-processing inequality is $\mathbb{H}(X|Y) \leq \mathbb{H}(X|Z)$, which means that, on average, processing Y can only increase our uncertainty about X .

Lemma 2.10 (Fano's inequality). *Let $X \in \mathcal{X}$ be a discrete random variable and let X' be any estimate of X that takes values in the same alphabet \mathcal{X} . Let $P_e \triangleq \mathbb{P}[X \neq X']$ be the probability of error obtained when estimating X with X' . Then,*

$$\mathbb{H}(X|X') \leq \mathbb{H}_b(P_e) + P_e \log(|\mathcal{X}| - 1),$$

where $\mathbb{H}_b(P_e)$ is the binary entropy function defined earlier.

Fano's inequality is the key ingredient of many proofs in this book because it relates an information-theoretic quantity (the conditional entropy $\mathbb{H}(X|X')$) to an operational quantity (the probability of error P_e). In what follows, we often write Fano's inequality in the form $\mathbb{H}(X|X') \leq \delta(P_e)$ to emphasize that $\mathbb{H}(X|X')$ goes to zero if P_e goes to zero.

Definition 2.6. *A function $f : \mathcal{I} \rightarrow \mathbb{R}$ defined on a set \mathcal{I} is convex on \mathcal{I} if, for all $(x_1, x_2) \in \mathcal{I}^2$ and for all $\lambda \in [0, 1]$,*

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

If the equation above holds with strict inequality, f is strictly convex on \mathcal{I} . A function $f : \mathcal{I} \rightarrow \mathbb{R}$ defined on a set \mathcal{I} is (strictly) concave on \mathcal{I} if the function $-f$ is (strictly) convex on \mathcal{I} .

Lemma 2.11 (Jensen's inequality). *Let $X \in \mathcal{X}$ be a random variable and let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a convex function. Then,*

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X]).$$

If f is strictly convex, then equality holds if and only if X is a constant.

Lemma 2.12. *Let $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$ be two discrete random variables with joint distribution p_{XY} . Then,*

- $\mathbb{H}(X)$ is a concave function of p_X ;
- $\mathbb{I}(X; Y)$ is a concave function of p_X for $p_{Y|X}$ fixed;
- $\mathbb{I}(X; Y)$ is a convex function of $p_{Y|X}$ for p_X fixed.

If X is a continuous random variable then, in general, $\mathbb{H}(X)$ is not well defined. It is convenient to define the differential entropy as follows.

Definition 2.7. Let $X \in \mathcal{X}$ be a continuous random variable with distribution p_X . The differential entropy of X is defined as

$$h(X) \triangleq - \int_{x \in \mathcal{X}} p_X(x) \log p_X(x) dx.$$

The notions of joint entropy, conditional entropy, and mutual information for continuous random variables are identical to their discrete counterparts, but with the differential entropy h in place of the entropy H . With the exception of Lemma 2.5 and Fano's inequality, all of the properties of entropy and mutual information stated above hold also for continuous random variables. In addition, the following properties will be useful in the next chapters.

Lemma 2.13. If $X \in \mathcal{X}$ is a continuous random variable with variance $\text{Var}(X) \leq \sigma^2$, then

$$h(X) \leq \frac{1}{2} \log(2\pi e \sigma^2),$$

with equality if and only if X has a Gaussian distribution with variance σ^2 .

Lemma 2.14 (Entropy–power inequality). Let $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$ be independent continuous random variables with entropy $h(X)$ and $h(Y)$, respectively. Let X' and Y' be independent Gaussian random variables such that $h(X') = h(X)$ and $h(Y') = h(Y)$. Then,

$$h(X + Y) \geq h(X' + Y'),$$

or, equivalently,

$$2^{2h(X+Y)} \geq 2^{2h(X')} + 2^{2h(Y')}.$$

2.1.3 Strongly typical sequences

Let $x^n \in \mathcal{X}^n$ be a sequence whose n elements are in a finite alphabet \mathcal{X} . The number of occurrences of a symbol $a \in \mathcal{X}$ in the sequence x^n is denoted by $N(a; x^n)$, and the empirical distribution (or histogram) of x^n is defined as the set $\{N(a; x^n)/n : a \in \mathcal{X}\}$.

Definition 2.8 (Strong typical set). Let p_X be a distribution on a finite alphabet \mathcal{X} and let $\epsilon > 0$. A sequence $x^n \in \mathcal{X}^n$ is (strongly) ϵ -typical with respect to p_X if

$$\forall a \in \mathcal{X} \quad \left| \frac{1}{n} N(a; x^n) - p_X(a) \right| \leq \epsilon p_X(a).$$

The set of all ϵ -typical sequences with respect to p_X is called the strong typical set and is denoted by $T_\epsilon^n(X)$.

In other words, the typical set $T_\epsilon^n(X)$ contains all sequences x^n whose empirical distribution is “close” to p_X . The notion of typicality is particularly useful in information theory because of a result known as the asymptotic equipartition property (AEP).

Theorem 2.1 (AEP). Let p_X be a distribution on a finite alphabet \mathcal{X} and let $0 < \epsilon < \min_{x \in \mathcal{X}} p_X(x)$. Let X^n be a sequence of independent and identically distributed (i.i.d.) random variables with distribution p_X . Then,

$$\begin{aligned} 1 - \delta_\epsilon(n) &\leq \mathbb{P}[X^n \in \mathcal{T}_\epsilon^n(X)] \leq 1, \\ (1 - \delta_\epsilon(n))2^{n(\mathbb{H}(X) - \delta(\epsilon))} &\leq |\mathcal{T}_\epsilon^n(X)| \leq 2^{n(\mathbb{H}(X) + \delta(\epsilon))}, \\ \forall x^n \in \mathcal{T}_\epsilon^n(X) \quad 2^{-n(\mathbb{H}(X) + \delta(\epsilon))} &\leq p_{X^n}(x^n) \leq 2^{-n(\mathbb{H}(X) - \delta(\epsilon))}. \end{aligned}$$

In simple terms, the AEP states that, for sufficiently large n , the probability that the realization x^n of a sequence of i.i.d. random variables belongs to the typical set is close to unity. Moreover, for practical purposes we may assume that the probability of any strongly typical sequence is about $2^{-n\mathbb{H}(X)}$ and the number of strongly typical sequences is approximately $2^{n\mathbb{H}(X)}$. In some sense, the AEP provides an operational interpretation of entropy.

Remark 2.1. It is possible to provide explicit expressions for $\delta_\epsilon(n)$ and $\delta(\epsilon)$ [6], but the rough characterization used in Theorem 2.1 is sufficient for our purposes. In particular, it makes it easier to keep track of small terms that depend on ϵ or n because we can write equations such as $\delta(\epsilon) + \delta(\epsilon) = \delta(\epsilon)$ without worrying about the exact dependence on ϵ .

The notion of typicality generalizes to multiple random variables. Assume that $(x^n, y^n) \in \mathcal{X}^n \times \mathcal{Y}^n$ is a pair of sequences with elements in finite alphabets \mathcal{X} and \mathcal{Y} . The number of occurrences of a pair $(a, b) \in \mathcal{X} \times \mathcal{Y}$ in the pair of sequences (x^n, y^n) is denoted by $N(a, b; x^n, y^n)$.

Definition 2.9 (Jointly typical set). Let p_{XY} be a joint distribution on the finite alphabets $\mathcal{X} \times \mathcal{Y}$ and let $\epsilon > 0$. Sequences $x^n \in \mathcal{X}^n$ and $y^n \in \mathcal{Y}^n$ are ϵ -jointly typical with respect to p_{XY} if

$$\forall (a, b) \in \mathcal{X} \times \mathcal{Y} \quad \left| \frac{1}{n} N(a, b; x^n, y^n) - p_{XY}(a, b) \right| \leq \epsilon p_{XY}(a, b).$$

The set of all ϵ -jointly typical sequences with respect to p_{XY} is called the jointly typical set and is denoted by $\mathcal{T}_\epsilon^n(XY)$.

One can check that $\mathcal{T}_\epsilon^n(XY) \subseteq \mathcal{T}_\epsilon^n(X) \times \mathcal{T}_\epsilon^n(Y)$. In other words, $(x^n, y^n) \in \mathcal{T}_\epsilon^n(XY)$ implies that $x^n \in \mathcal{T}_\epsilon^n(X)$ and $y^n \in \mathcal{T}_\epsilon^n(Y)$. This property is known as the *consistency* of joint typicality. Notice that the jointly typical set $\mathcal{T}_\epsilon^n(XY)$ is the typical set $\mathcal{T}_\epsilon^n(Z)$ for the random variable $Z = (X, Y)$. Therefore, the result below follows directly from Theorem 2.1.

Corollary 2.1 (Joint AEP). Let p_{XY} be a joint distribution on the finite alphabets $\mathcal{X} \times \mathcal{Y}$ and let $0 < \epsilon < \min_{(x, y) \in \mathcal{X} \times \mathcal{Y}} p_{XY}(x, y)$. Let (X^n, Y^n) be a sequence of i.i.d. random variables with joint distribution p_{XY} . Then,

$$\begin{aligned} 1 - \delta_\epsilon(n) &\leq \mathbb{P}[(X^n, Y^n) \in \mathcal{T}_\epsilon^n(XY)] \leq 1, \\ (1 - \delta_\epsilon(n))2^{n(\mathbb{H}(XY) - \delta(\epsilon))} &\leq |\mathcal{T}_\epsilon^n(XY)| \leq 2^{n(\mathbb{H}(XY) + \delta(\epsilon))}, \\ \forall (x^n, y^n) \in \mathcal{T}_\epsilon^n(XY) \quad 2^{-n(\mathbb{H}(XY) + \delta(\epsilon))} &\leq p_{X^n Y^n}(x^n, y^n) \leq 2^{-n(\mathbb{H}(XY) - \delta(\epsilon))}. \end{aligned}$$

It is also useful to introduce a conditional typical set, for which we can establish a conditional version of the AEP.

Definition 2.10. Let p_{XY} be a joint distribution on the finite alphabets $\mathcal{X} \times \mathcal{Y}$ and let $\epsilon > 0$. Let $x^n \in \mathcal{T}_\epsilon^n(X)$. The set

$$\mathcal{T}_\epsilon^n(XY|x^n) \triangleq \{y^n \in \mathcal{Y}^n : (x^n, y^n) \in \mathcal{T}_\epsilon^n(XY)\}$$

is called the conditional typical set with respect to x^n .

Theorem 2.2 (Conditional AEP). Let p_{XY} be a joint distribution on the finite alphabets $\mathcal{X} \times \mathcal{Y}$ and suppose that $0 < \epsilon' < \epsilon \leq \min_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p_{XY}(x, y)$. Let $x^n \in \mathcal{T}_{\epsilon'}^n(X)$ and let \tilde{Y}^n be a sequence of random variables such that

$$\forall y^n \in \mathcal{Y}^n \quad p_{\tilde{Y}^n}(y^n) = \prod_{i=1}^n p_{Y|X}(y_i|x_i).$$

Then,

$$\begin{aligned} 1 - \delta_{\epsilon\epsilon'}(n) &\leq \mathbb{P}[\tilde{Y}^n \in \mathcal{T}_\epsilon^n(XY|x^n)] \leq 1, \\ (1 - \delta_{\epsilon\epsilon'}(n))2^{n(\mathbb{H}(Y|X) - \delta(\epsilon))} &\leq |\mathcal{T}_\epsilon^n(XY|x^n)| \leq 2^{n(\mathbb{H}(Y|X) + \delta(\epsilon))}, \\ \forall y^n \in \mathcal{T}_\epsilon^n(XY|x^n) \quad 2^{-n(\mathbb{H}(Y|X) + \delta(\epsilon))} &\leq p_{Y^n|X^n}(y^n|x^n) \leq 2^{-n(\mathbb{H}(Y|X) - \delta(\epsilon))}. \end{aligned}$$

The conditional AEP means that, if x^n is a typical sequence and if \tilde{Y}^n is distributed according to $\prod_{i=1}^n p_{Y|X}(y_i|x_i)$, then \tilde{Y}^n is jointly typical with x^n with high probability for n large enough. In addition, the number of sequences y^n that are jointly typical with x^n is approximately $2^{n\mathbb{H}(Y|X)}$, and their probability is on the order of $2^{-n\mathbb{H}(Y|X)}$. The following corollary of the conditional AEP will be useful.

Corollary 2.2. Let p_{XY} be a joint distribution on the finite alphabets $\mathcal{X} \times \mathcal{Y}$ and let $0 < \epsilon < \mu_{XY}$ with $\mu_{XY} \triangleq \min_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p_{XY}(x, y)$. Let \tilde{Y}^n be a sequence of i.i.d. random variables with distribution p_Y . Then,

- if $x^n \in \mathcal{T}_\epsilon^n(X)$,

$$(1 - \delta_\epsilon(n))2^{-n(\mathbb{I}(X;Y) + \delta(\epsilon))} \leq \mathbb{P}[\tilde{Y}^n \in \mathcal{T}_\epsilon^n(XY|x^n)] \leq 2^{-n(\mathbb{I}(X;Y) - \delta(\epsilon))};$$

- if \tilde{X}^n is a sequence of random variables independent of \tilde{Y}^n and with arbitrary distribution $p_{\tilde{X}^n}$ on \mathcal{X}^n ,

$$\mathbb{P}[(\tilde{X}^n, \tilde{Y}^n) \in \mathcal{T}_\epsilon^n(XY)] \leq 2^{-n(\mathbb{I}(X;Y) - \delta(\epsilon))}.$$

In other words, if \tilde{Y}^n is generated independently of x^n , the probability that \tilde{Y}^n is jointly typical with x^n is small and on the order of $2^{-n\mathbb{I}(X;Y)}$. Corollary 2.2 generalizes to more than two random variables; in particular, we make extensive use of the following result.

Corollary 2.3. Let p_{UXY} be a joint distribution on the finite alphabets $\mathcal{U} \times \mathcal{X} \times \mathcal{Y}$ and suppose $0 < \epsilon < \mu_{UXY}$ with $\mu_{UXY} \triangleq \min_{(u,x,y) \in \mathcal{U} \times \mathcal{X} \times \mathcal{Y}} p_{UXY}(u, x, y)$. Let $(\tilde{U}^n, \tilde{X}^n)$ be

a sequence of random variables with arbitrary distribution $p_{\tilde{U}^n \tilde{X}^n}$ on $\mathcal{U}^n \times \mathcal{X}^n$. Let \tilde{Y}^n be a sequence of random variables conditionally independent of \tilde{X}^n given \tilde{U}^n such that

$$\forall (u^n, x^n, y^n) \in \mathcal{U}^n \times \mathcal{X}^n \times \mathcal{Y}^n$$

$$p_{\tilde{U}^n \tilde{X}^n \tilde{Y}^n}(u^n, x^n, y^n) = \left(\prod_{i=1}^n p_{Y|U}(y_i | u_i) \right) p_{\tilde{U}^n \tilde{X}^n}(u^n, x^n).$$

Then,

$$\mathbb{P}_{\tilde{U}^n \tilde{X}^n \tilde{Y}^n}[(\tilde{U}^n, \tilde{X}^n, \tilde{Y}^n) \in \mathcal{T}_\epsilon^n(\mathcal{U}\mathcal{X}\mathcal{Y})] \leq 2^{-n(\mathbb{H}(\mathcal{X}; \mathcal{Y}|\mathcal{U}) - \delta(\epsilon))}.$$

2.1.4 Weakly typical sequences

Strong typicality requires the relative frequency of each possible symbol to be close to the corresponding probability; however, the notion of strong typicality does not apply to continuous random variables and it is sometimes convenient to use a *weaker* notion of typicality, which merely requires the empirical entropy of a sequence to be close to the true entropy of the corresponding random variable. All definitions and results in this section are stated for discrete random variables but hold also for continuous random variables on replacing the entropy \mathbb{H} with the differential entropy \mathbb{h} .

Definition 2.11 (Weakly typical set). Let p_X be a distribution on a finite alphabet \mathcal{X} and let $\epsilon > 0$. A sequence $x^n \in \mathcal{X}^n$ is (weakly) ϵ -typical with respect to p_X if

$$\left| -\frac{1}{n} \log p_{X^n}(x^n) - \mathbb{H}(X) \right| \leq \epsilon.$$

The set of all weakly ϵ -typical sequences with respect to p_X is called the weakly typical set and is denoted $\mathcal{A}_\epsilon^n(X)$.

The weak version of the AEP then follows from the weak law of large numbers.

Theorem 2.3 (AEP). Let p_X be a distribution on a finite alphabet \mathcal{X} and let $\epsilon > 0$. Let X^n be a sequence of i.i.d. random variables with distribution p_X . Then,

- for n sufficiently large, $\mathbb{P}[X^n \in \mathcal{A}_\epsilon^n(X)] > 1 - \epsilon$;
- if $x^n \in \mathcal{A}_\epsilon^n(X)$, then $2^{-n(\mathbb{H}(X)+\epsilon)} \leq p_{X^n}(x^n) \leq 2^{-n(\mathbb{H}(X)-\epsilon)}$;
- for n sufficiently large, $(1 - \epsilon)2^{n(\mathbb{H}(X)-\epsilon)} \leq |\mathcal{A}_\epsilon^n(X)| \leq 2^{n(\mathbb{H}(X)+\epsilon)}$.

Definition 2.12 (Jointly weakly typical set). Let $p_{X\mathcal{Y}}$ be a joint distribution on the finite alphabets $\mathcal{X} \times \mathcal{Y}$ and let $\epsilon > 0$. Sequences $x^n \in \mathcal{X}^n$ and $y^n \in \mathcal{Y}^n$ are jointly (weakly) ϵ -typical with respect to $p_{X\mathcal{Y}}$ if

$$\begin{aligned} \left| -\frac{1}{n} \log p_{X^n Y^n}(x^n, y^n) - \mathbb{H}(XY) \right| &\leq \epsilon, \\ \left| -\frac{1}{n} \log p_{X^n}(x^n) - \mathbb{H}(X) \right| &\leq \epsilon, \\ \left| -\frac{1}{n} \log p_{Y^n}(y^n) - \mathbb{H}(Y) \right| &\leq \epsilon. \end{aligned}$$

The set of all jointly weakly ϵ -typical sequences with respect to p_{XY} is called the jointly weakly typical set and is denoted $\mathcal{A}_\epsilon^n(XY)$.

Theorem 2.4 (joint AEP). Let p_{XY} be a joint distribution on the finite alphabets $\mathcal{X} \times \mathcal{Y}$ and let $\epsilon > 0$. Let (X^n, Y^n) be a sequence of i.i.d. random variables with joint distribution p_{XY} . Then,

- for n sufficiently large, $\mathbb{P}[(X^n, Y^n) \in \mathcal{A}_\epsilon^n(XY)] > 1 - \epsilon$;
- if $(x^n, y^n) \in \mathcal{A}_\epsilon^n(XY)$, then $2^{-n(\mathbb{H}(XY)+\epsilon)} \leq p_{X^n Y^n}(x^n, y^n) \leq 2^{-n(\mathbb{H}(XY)-\epsilon)}$;
- for n sufficiently large, $(1 - \epsilon)2^{n(\mathbb{H}(XY)-\epsilon)} \leq |\mathcal{A}_\epsilon^n(XY)| \leq 2^{n(\mathbb{H}(XY)+\epsilon)}$.

With weak typicality, there is no exact counterpart to the conditional AEP given in Corollary 2.2 but the following result holds nevertheless.

Theorem 2.5. Let p_{XY} be a joint distribution on the finite alphabets $\mathcal{X} \times \mathcal{Y}$ and let $\epsilon > 0$. Let \tilde{Y}^n be a sequence of i.i.d. random variables with distribution p_Y , and let \tilde{X}^n be an independent sequence of i.i.d. random variables with distribution p_X . Then,

$$\mathbb{P}[(\tilde{X}^n, \tilde{Y}^n) \in \mathcal{A}_\epsilon^n(XY)] \leq 2^{-n(\mathbb{I}(X;Y)-\delta(\epsilon))}.$$

In subsequent chapters, we use the term AEP for both strong and weak typicality; however, it will be clear from the context whether we refer to the theorems of Section 2.1.3 or those of Section 2.1.4.

2.1.5 Markov chains and functional dependence graphs

The identification of Markov chains among random variables that depend on each other via complicated relations is a recurrent problem in information theory. In principle, Markov chains can be identified by manipulating the joint probability distribution of random variables, but this is often a tedious task. In this short section, we describe a graphical yet correct method for identifying Markov chains that is based on the *functional dependence graph* of random variables.

Definition 2.13 (Functional dependence graph). Consider m independent random variables and n functions of these variables. A functional dependence graph is a directed graph having $m + n$ vertices, and in which edges are drawn from one vertex to another if the random variable of the former vertex is an argument in the function defining the latter.

Example 2.1. Let $M \in \mathcal{M}$ and $Z^n \in \mathbb{R}^n$ be independent random variables. Let $\{f_i\}_n$ be a set of functions from \mathcal{M} to \mathbb{R}^n . For $i \in \llbracket 1, n \rrbracket$ define the random variables $X_i = f_i(M)$ and $Y_i = X_i + Z_i$. The functional dependence graph of the random variables M, X^n, Y^n , and Z^n is shown in Figure 2.1.

Definition 2.14 (d-separation). Let \mathcal{X} , \mathcal{Y} , and \mathcal{Z} be disjoint subsets of vertices in a functional dependence graph \mathcal{G} . The subset \mathcal{Z} is said to d-separate \mathcal{X} from \mathcal{Y} if there

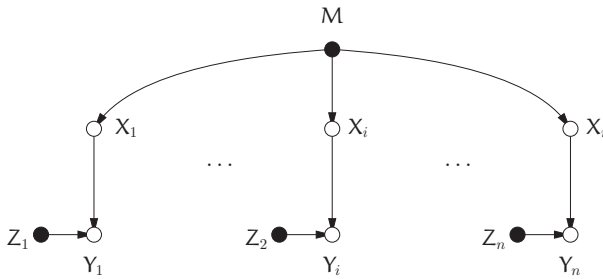


Figure 2.1 Functional dependence graph of variables in Example 2.1. For clarity, the independent random variables are indicated by filled circles (\bullet) whereas the functions of these random variables are indicated by empty circles (\circ).

exists no path between a vertex of \mathcal{X} and a vertex of \mathcal{Y} after the following operations have been performed:

- construct the subgraph \mathcal{G}' consisting of all vertices in \mathcal{X} , \mathcal{Y} , and \mathcal{Z} , as well as the edges and vertices encountered when moving backward starting from any of the vertices in \mathcal{X} , \mathcal{Y} , or \mathcal{Z} ;
- in the subgraph \mathcal{G}' , delete all edges coming out of \mathcal{Z} ;
- remove all arrows in \mathcal{G}' to obtain an undirected graph.

The usefulness of d-separation is justified by the following theorem.

Theorem 2.6. *Let \mathcal{X} , \mathcal{Y} , and \mathcal{Z} be disjoint subsets of the vertices in a functional dependence graph. If \mathcal{Z} d-separates \mathcal{X} from \mathcal{Y} , and if we collect the random variables in \mathcal{X} , \mathcal{Y} , and \mathcal{Z} in the random vectors X , Y , and Z , respectively, then $X \rightarrow Z \rightarrow Y$ forms a Markov chain.*

Theorem 2.6 is particularly useful in the converse proofs of channel coding theorems.

Example 2.2. On the basis of the functional dependence graph of Figure 2.1, one can check that, for any $i \neq j$, $X_i \rightarrow X_j \rightarrow Y_j$.

2.2 The point-to-point communication problem

The foundations of information theory were laid by Claude E. Shannon in his 1948 paper “A mathematical theory of communication” [7]. In his own words, *the fundamental problem of communication is that of reproducing at one point either exactly or approximately a message selected at another point*. If the message – for example, a letter from the alphabet, the gray level of a pixel or some physical quantity measured by a sensor – is to be reproduced at a remote location with a certain fidelity, some amount of information must be transmitted over a physical channel. This observation is the basis of Shannon’s general model for *point-to-point* communication reproduced in Figure 2.2. It consists of the following elements.

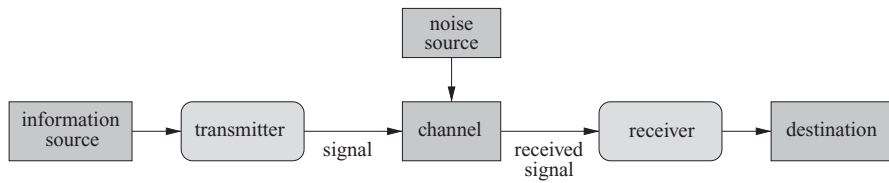


Figure 2.2 Shannon's communication model (from [7]).

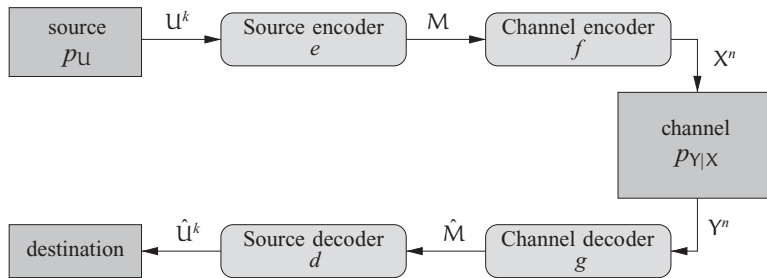


Figure 2.3 Mathematical model of a two-stage communication system.

- The *information source* generates messages according to some random process.
- The *transmitter* observes the messages and forms a signal to be sent over the channel.
- The *channel* is governed by a *noise source*, which corrupts the original input signal; this models the physical constraints of a communication system, for instance thermal noise in electronic circuits or multipath fading in a wireless medium.
- The *receiver* takes the received signal, forms a reconstructed version of the original message, and delivers the result to the *destination*.

Given the statistical properties of the information source and the noisy channel, the goal of the communication engineer is to design the transmitter and the receiver in a way that allows the information sent by the source to reach its destination in a *reliable* way. Information theory can help us achieve this goal by characterizing the fundamental mechanisms behind communication systems and providing us with precise mathematical conditions under which reliable communication is possible.

2.2.1 Point-to-point communication model

To give a precise formulation of the point-to-point communication problem, we require definitions for each of its constituent modules. We assume that the *source* and the *channel* are described by discrete-time random processes, and we determine that the receiver and the transmitter agree on a common *code*, specified by an *encoder* and *decoder* pair. As illustrated in Figure 2.3, we consider a two-stage system in which the source is compressed before being encoded for channel transmission, and channel outputs are decoded before being decompressed. The basic relationships among the components in Figure 2.3 are described in the following lines.

Definition 2.15. A discrete memoryless source (DMS) $(\mathcal{U}, p_{\mathcal{U}})$ generates a sequence of i.i.d. symbols (or letters) from the finite alphabet \mathcal{U} according to the probability distribution $p_{\mathcal{U}}$. The random variable representing a source symbol is denoted by \mathcal{U} .

Definition 2.16. A discrete memoryless channel (DMC) $(\mathcal{X}, p_{Y|X}, \mathcal{Y})$ is described by a finite input alphabet \mathcal{X} , a finite output alphabet \mathcal{Y} , and a conditional probability distribution $p_{Y|X}$, such that X and Y denote the channel input and the channel output, respectively. The set of conditional probabilities (also called transition probabilities) can be represented by a channel transition probability matrix $(p_{Y|X}(y|x))_{x,y}$.

In what follows, we illustrate many results numerically with the following DMCs.

Example 2.3. A binary symmetric channel with cross-over probability $p \in [0, 1]$, denoted by $\text{BSC}(p)$, is a DMC $(\{0, 1\}, p_{Y|X}, \{0, 1\})$ characterized by the transition probability matrix

$$\begin{pmatrix} 1-p & p \\ p & 1-p \end{pmatrix}.$$

Example 2.4. A binary erasure channel with erasure probability $\epsilon \in [0, 1]$, denoted by $\text{BEC}(\epsilon)$, is a DMC $(\{0, 1\}, p_{Y|X}, \{0, ?, 1\})$ characterized by the transition probability matrix

$$\begin{pmatrix} 1-\epsilon & \epsilon & 0 \\ 0 & \epsilon & 1-\epsilon \end{pmatrix}.$$

Definition 2.17. A $(2^{kR}, k)$ source code \mathcal{C}_k for a DMS $(\mathcal{U}, p_{\mathcal{U}})$ consists of

- a message set $\mathcal{M} = \llbracket 1, 2^{kR} \rrbracket$;
- an encoding function $e : \mathcal{U}^k \rightarrow \mathcal{M}$, which maps a sequence of k source symbols u^k to a message m ;
- a decoding function $d : \mathcal{M} \rightarrow \mathcal{U}^k \cup \{?\}$, which maps a message m to a sequence of source symbols $\hat{u}^k \in \mathcal{U}^k$ or an error message $?$.

The compression rate of the source code is defined as $(1/k)\log\lceil 2^{kR} \rceil$ in bits¹ per source symbol, and its probability of error is

$$\mathbf{P}_e(\mathcal{C}_k) \triangleq \mathbb{P}[\hat{\mathcal{U}}^k \neq \mathcal{U}^k \mid \mathcal{C}_k].$$

Definition 2.18. A rate R is an achievable compression rate for the source $(\mathcal{U}, p_{\mathcal{U}})$ if there exists a sequence of $(2^{kR}, k)$ source codes $\{\mathcal{C}_k\}_{k \geq 1}$, such that

$$\lim_{k \rightarrow \infty} \mathbf{P}_e(\mathcal{C}_k) = 0,$$

that is, the source sequences can be reconstructed with arbitrarily small probability of error with compression rates arbitrarily close to R .

¹ Unless specified otherwise, all logarithms are taken to the base two.

Definition 2.19. A $(2^{nR}, n)$ channel code C_n for a DMC $(\mathcal{X}, p_{Y|X}, \mathcal{Y})$ consists of

- a message set $\mathcal{M} = \llbracket 1, 2^{nR} \rrbracket$;
- an encoding function $f : \mathcal{M} \rightarrow \mathcal{X}^n$, which maps a message m to a codeword x^n with n symbols;
- a decoding function $g : \mathcal{Y}^n \rightarrow \mathcal{M} \cup \{?\}$, which maps a block of n channel outputs y^n to a message $\hat{m} \in \mathcal{M}$ or an error message $?$.

The set of codewords $\{f(m) : m \in \llbracket 1, 2^{nR} \rrbracket\}$ is called the *codebook* of C_n . With a slight abuse of notation, we denote the codebook itself by C_n as well. Unless specified otherwise, messages are represented by a random variable M uniformly distributed in \mathcal{M} , and the rate of the channel code is defined as $(1/n) \log \lceil 2^{nR} \rceil$ in bits per channel use. The average probability of error is defined as

$$\mathbf{P}_e(C_n) \triangleq \mathbb{P}[\hat{M} \neq M \mid C_n].$$

Definition 2.20. A rate R is an *achievable transmission rate* for the DMC $(\mathcal{X}, p_{Y|X}, \mathcal{Y})$ if there exists a sequence of $(2^{nR}, n)$ codes $\{C_n\}_{n \geq 1}$ such that

$$\lim_{n \rightarrow \infty} \mathbf{P}_e(C_n) = 0;$$

that is, messages can be transmitted at a rate arbitrarily close to R and decoded with arbitrarily small probability of error. The *channel capacity* of the DMC is defined as

$$C \triangleq \sup\{R : R \text{ is an achievable transmission rate}\}.$$

The typical goal of information theory is to characterize achievable rates on the basis of information-theoretic quantities that depend only on the given probability distributions and not on the block lengths k or n . A theorem that confirms the existence of codes for a class of achievable rates is often referred to as a *direct result* and the arguments that lead to this result constitute an *achievability proof*. On the other hand, when a theorem asserts that codes with certain properties do not exist, we speak of a *converse result* and a *converse proof*. A fundamental result that includes both the achievability and the converse parts is called a *coding theorem*. The mathematical tools that enable this characterization are those presented in Section 2.1, and we illustrate their use by discussing two of Shannon's fundamental coding theorems. These results form the basis of information theory and are of great use in several of the proofs developed in subsequent chapters.

Remark 2.2. Notice that the formulation of the point-to-point communication problem does not put any constraints either on the computational complexity or on the delay of the encoding and decoding procedures. In other words, the goal is to describe the fundamental limits of communications systems irrespective of their technological limitations.

2.2.2 The source coding theorem

The *source coding theorem* gives a complete solution (achievability and converse) for the point-to-point communication problem stated in Section 2.2.1 when the channel

is noiseless, that is $Y = X$. In that case, it is not necessary to use a channel code to compensate for the impairments caused by the channel, but it is still useful to encode the messages produced by the source to achieve a more efficient representation of the source information in bits per source symbol. This procedure is called *source coding* or *data compression*. The main idea is to consider only a subset \mathcal{A} of all possible source sequences \mathcal{U}^k , and assign a different index $i \in \llbracket 1, |\mathcal{A}| \rrbracket$ to each of the sequences $u^k \in \mathcal{A}$. If the source produces a sequence $u^k \in \mathcal{A}$, then the encoder outputs the corresponding index i , otherwise it outputs some predefined constant. The decoder receives the index i and outputs the corresponding sequence in \mathcal{A} .

Since information theory is primarily concerned with the fundamental limits of reliable communication, it is possible to prove the existence of codes without having to search for explicit code constructions. One technique, which is particularly useful in information-theoretic problems related to source coding, consists of *throwing* sequences $u^n \in \mathcal{U}^n$ randomly into a finite set of bins, such that the sequences that land in the same bin share a common bin index. If each sequence is assigned a bin at random according to a uniform distribution, then we refer to this procedure as *random binning*. If we want to prove that there exists a code such that the error probability goes to zero, it suffices to show that the average of the probability of error taken over all possible bin assignments goes to zero and to use the selection lemma. The following theorem exploits random binning to characterize the set of achievable compression rates.

Theorem 2.7 (Source coding theorem). *For a discrete memoryless source $(\mathcal{U}, p_{\mathcal{U}})$,*

$$\inf\{R : R \text{ is an achievable compression rate}\} = \mathbb{H}(\mathcal{U}).$$

In other words, if a compression rate R satisfies $R > \mathbb{H}(\mathcal{U})$ then R is achievable and any achievable compression rate R must satisfy $R \geq \mathbb{H}(\mathcal{U})$.

Proof. We start with the achievability part of the proof, which is based on random binning. The idea is to randomly assign each source sequence to one of a finite number of bins; then, as long as the number of bins is larger than $2^{k\mathbb{H}(\mathcal{U})}$, the probability of finding more than one typical sequence in the same bin is very small. If each typical sequence is mapped to a different bin index, an arbitrarily small probability of error can be achieved by letting the decoder output the typical sequence that corresponds to the received index. Formally, let $\epsilon > 0$ and $k \in \mathbb{N}^*$. Let $R > 0$ be a rate to be specified later. We construct a $(2^{kR}, k)$ source code \mathcal{C}_k as follows.

- *Binning.* For each sequence $u^k \in \mathcal{T}_{\epsilon}^k(\mathcal{U})$, draw an index uniformly at random in the set $\llbracket 1, 2^{kR} \rrbracket$. The index assignment defines the encoding function

$$e : \mathcal{U}^k \rightarrow \llbracket 1, 2^{kR} \rrbracket,$$

which is revealed to the encoder and decoder.

- *Encoder.* Given an observation u^k , output $m = e(u^k)$ if $u^k \in \mathcal{T}_{\epsilon}^k(\mathcal{U})$; otherwise output $m = 1$.
- *Decoder.* Given message m , output \hat{u}^k if it is the unique sequence such that $\hat{u}^k \in \mathcal{T}_{\epsilon}^k(\mathcal{U})$ and $e(\hat{u}^k) = m$; otherwise output an error ?.

The random variable that represents the randomly generated encoding function e is denoted by \mathbb{E} while the random variable that represents the randomly generated code \mathcal{C}_k is denoted by \mathcal{C}_k . We proceed to bound $\mathbb{E}[\mathbf{P}_e(\mathcal{C}_k)]$. First, note that $\mathbb{E}[\mathbf{P}_e(\mathcal{C}_k)]$ can be expressed in terms of the events

$$\begin{aligned}\mathcal{E}_0 &= \{\mathbb{U}^k \notin \mathcal{T}_\epsilon^k(\mathbb{U})\}, \\ \mathcal{E}_1 &= \{\exists \hat{u}^k \neq \mathbb{U}^k : \mathbb{E}(\hat{u}^k) = \mathbb{E}(\mathbb{U}^k) \text{ and } \hat{u}^k \in \mathcal{T}_\epsilon^k(\mathbb{U})\}\end{aligned}$$

as $\mathbb{E}[\mathbf{P}_e(\mathcal{C}_k)] = \mathbb{P}[\mathcal{E}_0 \cup \mathcal{E}_1]$. By the union bound,

$$\mathbb{E}[\mathbf{P}_e(\mathcal{C}_k)] \leq \mathbb{P}[\mathcal{E}_0] + \mathbb{P}[\mathcal{E}_1]. \quad (2.1)$$

By the AEP,

$$\mathbb{P}[\mathcal{E}_0] \leq \delta_\epsilon(k) \quad (2.2)$$

and we can upper bound $\mathbb{P}[\mathcal{E}_1]$ as

$$\begin{aligned}\mathbb{P}[\mathcal{E}_1] &= \sum_{u^k} p_{\mathbb{U}^k}(u^k) \mathbb{P}[\exists \hat{u}^k \neq u^k : \mathbb{E}(\hat{u}^k) = \mathbb{E}(u^k) \text{ and } \hat{u}^k \in \mathcal{T}_\epsilon^k(\mathbb{U})] \\ &\leq \sum_{u^k} p_{\mathbb{U}^k}(u^k) \sum_{\substack{\hat{u}^k \in \mathcal{T}_\epsilon^k(\mathbb{U}) \\ \hat{u}^k \neq u^k}} \mathbb{P}[\mathbb{E}(\hat{u}^k) = \mathbb{E}(u^k)] \\ &= \sum_{u^k} p_{\mathbb{U}^k}(u^k) \sum_{\substack{\hat{u}^k \in \mathcal{T}_\epsilon^k(\mathbb{U}) \\ \hat{u}^k \neq u^k}} \frac{1}{\lceil 2^{kR} \rceil} \\ &\leq \sum_{u^k} p_{\mathbb{U}^k}(u^k) \frac{1}{\lceil 2^{kR} \rceil} |\mathcal{T}_\epsilon^k(\mathbb{U})| \\ &\leq \sum_{u^k} p_{\mathbb{U}^k}(u^k) \frac{1}{\lceil 2^{kR} \rceil} 2^{k(\mathbb{H}(\mathbb{U}) + \delta(\epsilon))} \\ &\leq 2^{k(\mathbb{H}(\mathbb{U}) + \delta(\epsilon) - R)}.\end{aligned}$$

Hence, if we choose $R > \mathbb{H}(\mathbb{U}) + \delta(\epsilon)$, we have

$$\mathbb{P}[\mathcal{E}_1] \leq \delta_\epsilon(k). \quad (2.3)$$

On substituting (2.2) and (2.3) into (2.1), we obtain $\mathbb{E}[\mathbf{P}_e(\mathcal{C}_k)] \leq \delta_\epsilon(k)$. By applying the selection lemma to the random variable \mathcal{C}_k and the function \mathbf{P}_e , we conclude that there exists at least one source code \mathcal{C}_k such that $\mathbf{P}_e(\mathcal{C}_k) \leq \delta_\epsilon(k)$. Since ϵ can be chosen arbitrarily small, all rates $R > \mathbb{H}(\mathbb{U})$ are achievable.

We now establish the converse result and show that any achievable rate must satisfy $R \geq \mathbb{H}(\mathbb{U})$. Let R be an achievable rate and let $\epsilon > 0$. By definition, there exists a source code \mathcal{C}_k such that $\mathbf{P}_e(\mathcal{C}_k) \leq \delta(\epsilon)$. If we let M denote the message output by the encoder,

then Fano's inequality guarantees that

$$\frac{1}{k} \mathbb{H}(\mathbf{U}^k | \mathcal{M}\mathcal{C}_k) \leq \delta(\mathbf{P}_\epsilon(\mathcal{C}_k)) \leq \delta(\epsilon).$$

We drop the conditioning on \mathcal{C}_k in subsequent calculations to simplify the notation. Note that

$$\begin{aligned} \mathbb{H}(\mathbf{U}) &= \frac{1}{k} \mathbb{H}(\mathbf{U}^k) \\ &= \frac{1}{k} \mathbb{I}(\mathbf{U}^k; \mathbf{M}) + \frac{1}{k} \mathbb{H}(\mathbf{U}^k | \mathbf{M}) \\ &\leq \frac{1}{k} \mathbb{I}(\mathbf{U}^k; \mathbf{M}) + \delta(\epsilon) \\ &\leq \frac{1}{k} \mathbb{H}(\mathbf{M}) + \delta(\epsilon) \\ &\leq R + \delta(k) + \delta(\epsilon). \end{aligned}$$

Since ϵ can be chosen arbitrarily small and k can be chosen arbitrarily large, we obtain $R \geq \mathbb{H}(\mathbf{U})$. \square

Remark 2.3. *Alternatively, the achievability part of the source coding theorem can be established on the basis of the AEP alone. In fact, for large k the AEP guarantees that any sequence \mathbf{u}^k produced by the source $(\mathcal{U}, p_{\mathbf{U}})$ belongs with high probability to the typical set $\mathcal{T}_\epsilon^k(\mathbf{U})$; hence, we need only index the approximately $2^{k\mathbb{H}(\mathbf{U})}$ typical sequences to achieve arbitrarily small probability of error and the corresponding rate is on the order of $\mathbb{H}(\mathbf{U})$.*

2.2.3 The channel coding theorem

The *channel coding theorem* gives a complete solution (achievability and converse) for the point-to-point communication problem stated in Section 2.2.1 when the source $(\mathcal{U}, p_{\mathbf{U}})$ is uniform over \mathcal{U} . According to the source coding theorem, there is no need to encode the source since $\mathbb{H}(\mathbf{U}) = \log|\mathcal{U}|$ is maximal. We simply group the source symbols in sequences of length k . Letting $M \triangleq |\mathcal{U}|^k$ and $\mathcal{M} = \llbracket 1, M \rrbracket$, we index each sequence of length k with an integer $m \in \mathcal{M}$. We use a channel code of rate $(1/n)\log M$ to transmit the messages produced by source \mathbf{U} over a discrete memoryless channel $(\mathcal{X}, p_{Y|X}, \mathcal{Y})$.

As was done for the source coding theorem, it is possible to prove the existence of codes without having to search for explicit code constructions. The idea is to construct a random code by drawing the symbols of codewords independently at random according to a fixed probability distribution p_X on \mathcal{X} . Then, if we want to prove that there exists a code such that the error probability goes to zero for n sufficiently large, it suffices to show that the average of the probability of error taken over all possible random codebooks goes to zero for n sufficiently large and use the selection lemma. This technique is

referred to as *random coding* and is used in the proof of the following theorem to characterize the set of achievable rates.

Theorem 2.8 (Channel coding theorem). *The capacity of a DMC $(\mathcal{X}, p_{Y|X}, \mathcal{Y})$ is $C = \max_{p_X} \mathbb{I}(X; Y)$. In other words, if $R < C$ then R is an achievable transmission rate and an achievable transmission rate must satisfy $R \leq C$.*

Proof. We begin with the achievability part based on random coding. We choose a probability distribution p_X on \mathcal{X} and, without loss of generality, we assume that p_X is such that $\mathbb{I}(X; Y) > 0$. Let $0 < \epsilon < \mu_{XY}$, where $\mu_{XY} = \min_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p_{XY}(x, y)$, and let $n \in \mathbb{N}^*$. Let $R > 0$ be a rate to be specified later. We construct a $(2^{nR}, n)$ code C_n as follows.

- *Codebook construction.* Construct a codebook with $\lceil 2^{nR} \rceil$ codewords, labeled $x^n(m)$ with $m \in \llbracket 1, 2^{nR} \rrbracket$, by generating the symbols $x_i(m)$ for $i \in \llbracket 1, n \rrbracket$ and $m \in \llbracket 1, 2^{nR} \rrbracket$ independently according to p_X . The codebook is revealed both to the encoder and to the decoder.
- *Encoder f .* Given m , transmit $x^n(m)$.
- *Decoder g .* Given y^n , output \hat{m} if it is the unique message such that $(x^n(\hat{m}), y^n) \in \mathcal{T}_\epsilon^n(XY)$; otherwise, output an error ?.

The random variable that represents the randomly generated codebook C_n is denoted by C_n . We first develop an upper bound for $\mathbb{E}[\mathbf{P}_e(C_n)]$. Notice that

$$\begin{aligned} \mathbb{E}[\mathbf{P}_e(C_n)] &= \mathbb{E}_{C_n} \left[\mathbb{P} \left[M \neq \hat{M} \mid C_n \right] \right] \\ &= \sum_m \mathbb{E}_{C_n} \left[\mathbb{P} \left[M \neq \hat{M} \mid M = m, C_n \right] \right] p_M(m). \end{aligned}$$

By virtue of the symmetry of the random code construction, we have that $\mathbb{E}_{C_n} [\mathbb{P} [M \neq \hat{M} \mid M = m, C_n]]$ is independent of m . Therefore, we can assume without losing generality that message $m = 1$ was sent and write

$$\mathbb{E}[\mathbf{P}_e(C_n)] = \mathbb{E}_{C_n} \left[\mathbb{P} \left[M \neq \hat{M} \mid M = 1, C_n \right] \right].$$

Notice that $\mathbb{E}[\mathbf{P}_e(C_n)]$ can be expressed in terms of the events

$$\mathcal{E}_i = \{(X^n(i), Y^n) \in \mathcal{T}_\epsilon^n(XY)\} \quad \text{for } i \in \llbracket 1, 2^{nR} \rrbracket$$

as $\mathbb{E}[\mathbf{P}_e(C_n)] = \mathbb{P} \left[\mathcal{E}_1^c \cup \bigcup_{i \neq 1} \mathcal{E}_i \right]$. By the union bound,

$$\mathbb{E}[\mathbf{P}_e(C_n)] \leq \mathbb{P} [\mathcal{E}_1^c] + \sum_{i \neq 1} \mathbb{P} [\mathcal{E}_i]. \quad (2.4)$$

By the AEP,

$$\mathbb{P} [\mathcal{E}_1^c] \leq \delta_\epsilon(n). \quad (2.5)$$

Since Y^n is the output of the channel when $X^n(1)$ is transmitted and since $X^n(1)$ is independent of $X^n(i)$ for $i \neq 1$, note that Y^n is independent of $X^n(i)$ for $i \neq 1$; hence,

Corollary 2.2 applies and

$$\mathbb{P}[\mathcal{E}_i] \leq 2^{-n(\mathbb{I}(X;Y) - \delta(\epsilon))} \quad \text{for } i \neq 1. \quad (2.6)$$

On substituting (2.5) and (2.6) into (2.4), we obtain

$$\mathbb{E}[\mathbf{P}_e(C_n)] \leq \delta_\epsilon(n) + \lceil 2^{nR} \rceil 2^{-n(\mathbb{I}(X;Y) - \delta(\epsilon))}.$$

Hence, if we choose the rate R such that $R < \mathbb{I}(X;Y) - \delta(\epsilon)$, then

$$\mathbb{E}[\mathbf{P}_e(C_n)] \leq \delta_\epsilon(n).$$

By applying the selection lemma to the random variable C_n and the function \mathbf{P}_e , we conclude that there exists a $(2^{nR}, n)$ code C_n such that $\mathbf{P}_e(C_n) \leq \delta_\epsilon(n)$. Since ϵ can be chosen arbitrarily small and since the distribution p_X is arbitrary, we conclude that all rates $R < \max_{p_X} \mathbb{I}(X;Y)$ are achievable.

We now establish the converse part of the proof. Let R be an achievable rate and let $\epsilon > 0$. For n sufficiently large, there exists a $(2^{nR}, n)$ code C_n such that

$$\frac{1}{n} \mathbb{H}(M|C_n) \geq R \quad \text{and} \quad \mathbf{P}_e(C_n) \leq \delta(\epsilon).$$

In the remaining part of the proof we drop the conditioning on C_n to simplify the notation. By virtue of Fano's inequality, it also holds that

$$\frac{1}{n} \mathbb{H}(M|Y^n) \leq \delta(\mathbf{P}_e(C_n)) = \delta(\epsilon).$$

Therefore,

$$\begin{aligned} R &\leq \frac{1}{n} \mathbb{H}(M) \\ &\leq \frac{1}{n} \mathbb{H}(M; Y^n) + \frac{1}{n} \mathbb{H}(M|Y^n) \\ &\leq \frac{1}{n} \mathbb{H}(M; Y^n) + \delta(\epsilon) \\ &\stackrel{(a)}{\leq} \frac{1}{n} \mathbb{H}(X^n; Y^n) + \delta(\epsilon) \\ &= \frac{1}{n} \mathbb{H}(Y^n) - \frac{1}{n} \mathbb{H}(Y^n|X^n) + \delta(\epsilon) \\ &\stackrel{(b)}{=} \frac{1}{n} \sum_{i=1}^n \left(\mathbb{H}(Y_i|Y^{i-1}) - \frac{1}{n} \mathbb{H}(Y_i|X_i) \right) + \delta(\epsilon) \\ &\stackrel{(c)}{\leq} \frac{1}{n} \sum_{i=1}^n \left(\mathbb{H}(Y_i) - \frac{1}{n} \mathbb{H}(Y_i|X_i) \right) + \delta(\epsilon) \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{H}(X_i; Y_i) + \delta(\epsilon) \\ &\leq \max_{p_X} \mathbb{H}(X; Y) + \delta(\epsilon), \end{aligned}$$

where (a) follows from the data-processing inequality applied to the Markov chain $M \rightarrow X^n \rightarrow Y^n$, (b) follows because the channel is memoryless, and (c) follows because conditioning does not increase entropy. Since ϵ can be chosen arbitrarily small, we obtain $R \leq \max_{p_X} \mathbb{I}(X; Y)$. \square

The channel coding theorem shows that the channel capacity is equal to the maximum mutual information between the channel input X and the channel output Y , where the maximization is carried out over all possible input probability distributions p_X . The proof technique and structure are common to most proofs in subsequent chapters.

Example 2.5. The capacity of a binary symmetric channel $\text{BSC}(p)$ is $1 - \mathbb{H}_b(p)$. The capacity of a binary erasure channel $\text{BEC}(\epsilon)$ is $1 - \epsilon$.

Among the many channel models, the additive white Gaussian noise (AWGN) channel (Gaussian channel for short) takes a particularly prominent role in information and communication theory, because it captures the impact of thermal noise and interference on wired and wireless communications. The channel output at each time $i \geq 1$ is given by $Y_i = X_i + N_i$, where X_i denotes the transmitted symbol and $\{N_i\}_{i \geq 1}$ are i.i.d. random variables with distribution $\mathcal{N}(0, \sigma^2)$. Since the channel capacity of the Gaussian channel can be infinite without further restrictions, we add an average power constraint in the form of

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i^2] \leq P.$$

Theorem 2.9. *The capacity of a Gaussian channel is given by*

$$C = \frac{1}{2} \log \left(1 + \frac{P}{\sigma^2} \right),$$

where P denotes the power constraint and σ^2 is the variance of the noise.

Sketch of proof. The proof developed for Theorem 2.8 does not apply directly to the Gaussian channel because of the power constraint imposed on channel inputs and the continuous nature of the channel. Nevertheless, it is possible to develop a similar proof by using weakly typical sequences and the weak AEP (see for instance [3, Chapter 9]). The power constraint can be dealt with by introducing an error event that accounts for the sequences violating the power constraint in the codebook generation. \square

2.3 Network information theory

Shannon's coding theorems characterize the fundamental limits of communication between two users. However, in many communication scenarios – for example, satellite broadcasting, cellular telephony, the Internet, and wireless sensor networks – the information is sent by one or more transmitting nodes to one or more receiving nodes

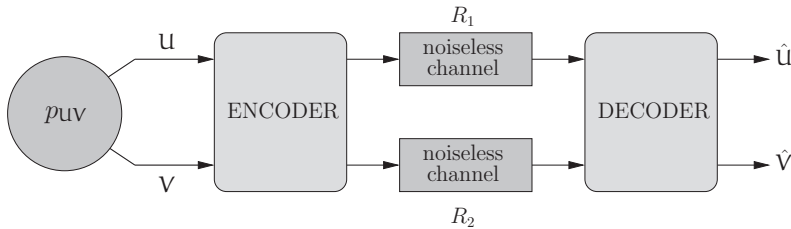


Figure 2.4 Joint encoding of correlated sources.

over more or less intricate communication networks. The interactions between the users of said networks introduce a whole new range of fundamental communication aspects that are not present in the classical point-to-point problem, such as *interference*, *user cooperation*, and *feedback*. The central goal of *network information theory* is to provide a thorough understanding of these basic mechanisms, by characterizing the fundamental limits of communication systems with multiple users. In this section, we discuss some results of network information theory that are useful for understanding information-theoretic security in subsequent chapters.

2.3.1 Distributed source coding

Consider a DMS (\mathcal{UV}, p_{UV}) that consists of two components \mathcal{U} and \mathcal{V} with joint distribution p_{UV} . As shown in Figure 2.4, the two components are to be processed by a joint encoder and transmitted to a common destination over two noiseless channels. The joint distribution p_{UV} can be arbitrary and the symbols produced by \mathcal{U} and \mathcal{V} at any given point in time are statistically dependent; therefore, we refer to \mathcal{U} and \mathcal{V} as *correlated sources*. Since the channels to the destination do not introduce any errors, we may ask the following question: at what rates R_1 and R_2 can we transmit information generated by \mathcal{U} and \mathcal{V} with an arbitrarily small probability of error? Since there is a common encoder and a common decoder, this problem reduces to the classical point-to-point problem and the solution follows naturally from Shannon's source coding theorem: the messages can be reconstructed with an arbitrarily small probability of error at the receiver if and only if

$$R_1 + R_2 > \mathbb{H}(\mathcal{UV});$$

that is, the sum rate must be greater than the joint entropy of \mathcal{U} and \mathcal{V} .

As illustrated in Figure 2.5, the problem becomes more challenging if instead of a joint encoder we consider two *separate* encoders. Here, each encoder observes only the realizations of the one source it is assigned to and does not know the output symbols of the other source. In this case, it is not immediately clear which encoding rates guarantee reconstruction with an arbitrarily small probability of error at the receiver. If we encode \mathcal{U} at rate $R_1 > \mathbb{H}(\mathcal{U})$ and \mathcal{V} at rate $R_2 > \mathbb{H}(\mathcal{V})$, then the source coding theorem guarantees once again that an arbitrarily small probability of error is possible. But, in this case, the sum rate satisfies $R_1 + R_2 > \mathbb{H}(\mathcal{U}) + \mathbb{H}(\mathcal{V})$, which, in general, is greater than the joint entropy $\mathbb{H}(\mathcal{UV})$.

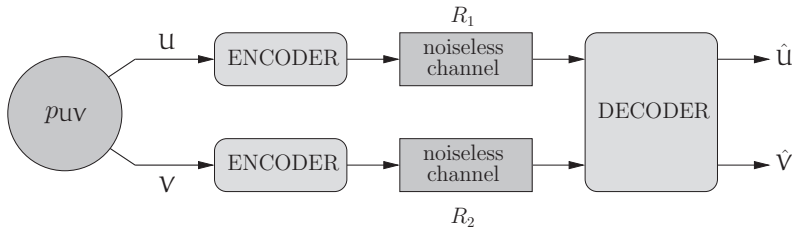


Figure 2.5 Separate encoding of correlated sources (the Slepian–Wolf problem).

Surprisingly, it turns out that the sum rate required by two separate encoders is the same as that required by a joint encoder, that is $R_1 + R_2 > \mathbb{H}(UV)$ is sufficient to reconstruct U and V with an arbitrarily small probability of error. In other words, there is no penalty in overall compression rate due to the fact that the encoders can observe only the realizations of the one source they have been assigned to. However, it is important to point out that the decoder does require a minimum amount of rate from each encoder; specifically, the average remaining uncertainty about the messages of one source given the messages of the other source, $\mathbb{H}(U|V)$ and $\mathbb{H}(V|U)$. Formally, a code for the distributed source coding problem is defined as follows.

Definition 2.21. A $(2^{kR_1}, 2^{kR_2}, k)$ source code \mathcal{C}_k for the DMS (UV, p_{UV}) consists of

- two message sets $\mathcal{M}_1 = [1, 2^{kR_1}]$ and $\mathcal{M}_2 = [1, 2^{kR_2}]$;
- an encoding function $e_1 : \mathcal{U}^k \rightarrow \mathcal{M}_1$, which maps a sequence of k source symbols u^k to a message m_1 ;
- an encoding function $e_2 : \mathcal{V}^k \rightarrow \mathcal{M}_2$, which maps a sequence of k source symbols v^k to a message m_2 ;
- a decoding function $d : \mathcal{M}_1 \times \mathcal{M}_2 \rightarrow (\mathcal{U}^k \times \mathcal{V}^k) \cup \{?\}$, which maps a message pair (m_1, m_2) to a pair of source sequences $(\hat{u}^k, \hat{v}^k) \in \mathcal{U}^k \times \mathcal{V}^k$ or an error message $?$.

The performance of a code \mathcal{C}_k is measured in terms of the average probability of error

$$\mathbf{P}_e(\mathcal{C}_k) \triangleq \mathbb{P}[(\hat{U}^k, \hat{V}^k) \neq (U^k, V^k) | \mathcal{C}_k].$$

Definition 2.22. A rate pair (R_1, R_2) is achievable if there exists a sequence of $(2^{kR_1}, 2^{kR_2}, k)$ codes $\{\mathcal{C}_k\}_{k \geq 1}$ such that

$$\lim_{k \rightarrow \infty} \mathbf{P}_e(\mathcal{C}_k) = 0.$$

The achievable rate region is defined as

$$\mathcal{R}^{\text{SW}} \triangleq \text{cl}(\{(R_1, R_2) : (R_1, R_2) \text{ is achievable}\}).$$

The achievable rate region with separate encoding was first characterized by Slepian and Wolf; hence, the region is often called the Slepian–Wolf region and codes for the distributed source coding problem are often referred to as Slepian–Wolf codes.

Theorem 2.10 (Slepian–Wolf theorem). *The achievable rate region with separate encoding for a source $(\mathcal{U}, \mathcal{V})$ is*

$$\mathcal{R}^{\text{SW}} \triangleq \left\{ (R_1, R_2): \begin{array}{l} R_1 \geq \mathbb{H}(\mathcal{U}|\mathcal{V}) \\ R_2 \geq \mathbb{H}(\mathcal{V}|\mathcal{U}) \\ R_1 + R_2 \geq \mathbb{H}(\mathcal{UV}) \end{array} \right\}.$$

Proof. We begin with the achievability part of the proof, which is based on joint typicality and random binning. Let $\epsilon > 0$ and $k \in \mathbb{N}^*$. Let $R_1 > 0$ and $R_2 > 0$ be rates to be specified later. We construct a $(2^{kR_1}, 2^{kR_2}, k)$ code \mathcal{C}_k as follows.

- *Binning.* For each sequence $u^k \in \mathcal{T}_\epsilon^k(\mathcal{U})$, draw an index uniformly at random in the set $[1, 2^{kR_1}]$. For each sequence $v^k \in \mathcal{T}_\epsilon^k(\mathcal{V})$, draw an index uniformly at random in the set $[1, 2^{kR_2}]$. The index assignments define the encoding functions

$$e_1 : \mathcal{U}^k \rightarrow [1, 2^{kR_1}] \quad \text{and} \quad e_2 : \mathcal{V}^k \rightarrow [1, 2^{kR_2}],$$

which are revealed to all parties.

- *Encoder 1.* Given the observation u^k , if $u^k \in \mathcal{T}_\epsilon^k(\mathcal{U})$, output $m_1 = e_1(u^k)$; otherwise output $m_1 = 1$.
- *Encoder 2.* Given the observation v^k , if $v^k \in \mathcal{T}_\epsilon^k(\mathcal{V})$, output $m_2 = e_2(v^k)$; otherwise output $m_2 = 1$.
- *Decoder.* Given messages m_1 and m_2 , output \hat{u}^k and \hat{v}^k if they are the unique sequences such that $(\hat{u}^k, \hat{v}^k) \in \mathcal{T}_\epsilon^k(\mathcal{UV})$ and $e_1(\hat{u}^k) = m_1, e_2(\hat{v}^k) = m_2$; otherwise, output ?.

The random variables that represent the randomly generated functions e_1 and e_2 are denoted by E_1 and E_2 , and the random variable that represents the randomly generated code \mathcal{C}_k is denoted by C_k . We proceed to bound $\mathbb{E}[\mathbf{P}_e(C_k)]$, which can be expressed in terms of the following events:

$$\begin{aligned} \mathcal{E}_0 &= \{(\mathcal{U}^k, \mathcal{V}^k) \notin \mathcal{T}_\epsilon^k(\mathcal{UV})\}, \\ \mathcal{E}_1 &= \{\exists \hat{u}^k \neq \mathcal{U}^k : E_1(\hat{u}^k) = E_1(\mathcal{U}^k) \text{ and } (\hat{u}^k, \mathcal{V}^k) \in \mathcal{T}_\epsilon^k(\mathcal{UV})\}, \\ \mathcal{E}_2 &= \{\exists \hat{v}^k \neq \mathcal{V}^k : E_2(\hat{v}^k) = E_2(\mathcal{V}^k) \text{ and } (\mathcal{U}^k, \hat{v}^k) \in \mathcal{T}_\epsilon^k(\mathcal{UV})\}, \\ \mathcal{E}_{12} &= \{\exists \hat{v}^k \neq \mathcal{V}^k, \hat{u}^k \neq \mathcal{U}^k, : E_1(\hat{u}^k) = E_1(\mathcal{U}^k), E_2(\hat{v}^k) = E_2(\mathcal{V}^k) \\ &\quad \text{and } (\hat{u}^k, \hat{v}^k) \in \mathcal{T}_\epsilon^k(\mathcal{UV})\}, \end{aligned}$$

since $\mathbb{E}[\mathbf{P}_e(C_k)] = \mathbb{P}[\mathcal{E}_0 \cup \mathcal{E}_1 \cup \mathcal{E}_2 \cup \mathcal{E}_{12}]$. By the union bound,

$$\mathbb{E}[\mathbf{P}_e(C_k)] \leq \mathbb{P}[\mathcal{E}_0] + \mathbb{P}[\mathcal{E}_1] + \mathbb{P}[\mathcal{E}_2] + \mathbb{P}[\mathcal{E}_{12}]. \quad (2.7)$$

By the AEP,

$$\mathbb{P}[\mathcal{E}_0] \leq \delta_\epsilon(k). \quad (2.8)$$

Using Theorem 2.2,

$$\begin{aligned}
 \mathbb{P}[\mathcal{E}_1] &= \sum_{u^k, v^k} p_{\mathcal{U}^k \mathcal{V}^k}(u^k, v^k) \mathbb{P}[\exists \hat{u}^k \neq u^k : E_1(\hat{u}^k) = E_1(u^k) \text{ and } (\hat{u}^k, v^k) \in \mathcal{T}_\epsilon^k(\mathcal{U}\mathcal{V})] \\
 &\leq \sum_{u^k, v^k} p_{\mathcal{U}^k \mathcal{V}^k}(u^k, v^k) \sum_{\substack{\hat{u}^k \in \mathcal{T}_\epsilon^k(\mathcal{U}\mathcal{V}|v^k) \\ \hat{u}^k \neq u^k}} \mathbb{P}[E_1(\hat{u}^k) = E_1(u^k)] \\
 &= \sum_{u^k, v^k} p_{\mathcal{U}^k \mathcal{V}^k}(u^k, v^k) \sum_{\substack{\hat{u}^k \in \mathcal{T}_\epsilon^k(\mathcal{U}\mathcal{V}|v^k) \\ \hat{u}^k \neq u^k}} \frac{1}{\lceil 2^{kR_1} \rceil} \\
 &\leq \sum_{u^k, v^k} p_{\mathcal{U}^k \mathcal{V}^k}(u^k, v^k) \frac{1}{\lceil 2^{kR_1} \rceil} |\mathcal{T}_\epsilon^k(\mathcal{U}\mathcal{V}|v^k)| \\
 &\leq \sum_{u^k, v^k} p_{\mathcal{U}^k \mathcal{V}^k}(u^k, v^k) \frac{1}{\lceil 2^{kR_1} \rceil} 2^{k(\mathbb{H}(\mathcal{U}|\mathcal{V}) + \delta(\epsilon))} \\
 &\leq 2^{k(\mathbb{H}(\mathcal{U}|\mathcal{V}) + \delta(\epsilon) - R_1)}. \tag{2.9}
 \end{aligned}$$

Similarly, we obtain the following bounds for $\mathbb{P}[\mathcal{E}_2]$ and $\mathbb{P}[\mathcal{E}_{12}]$:

$$\mathbb{P}[\mathcal{E}_2] \leq 2^{k(\mathbb{H}(\mathcal{V}|\mathcal{U}) + \delta(\epsilon) - R_2)}, \tag{2.10}$$

$$\mathbb{P}[\mathcal{E}_{12}] \leq 2^{k(\mathbb{H}(\mathcal{U}\mathcal{V}) + \delta(\epsilon) - (R_1 + R_2))}. \tag{2.11}$$

Hence, if we choose the rates R_1 and R_2 such that

$$R_1 > \mathbb{H}(\mathcal{U}|\mathcal{V}) + \delta(\epsilon),$$

$$R_2 > \mathbb{H}(\mathcal{V}|\mathcal{U}) + \delta(\epsilon),$$

$$R_1 + R_2 > \mathbb{H}(\mathcal{U}\mathcal{V}) + \delta(\epsilon),$$

and substitute (2.8)–(2.11) into (2.7), we obtain $\mathbb{E}[\mathbf{P}_\epsilon(C_k)] \leq \delta_\epsilon(k)$. By applying the selection lemma to the random variable C_k and the function \mathbf{P}_ϵ , we conclude that there exists a specific code C_k such that $\mathbf{P}_\epsilon(C_k) \leq \delta_\epsilon(k)$. Since ϵ can be chosen arbitrarily small, we conclude that

$$\left\{ (R_1, R_2): \begin{array}{l} R_1 \geq \mathbb{H}(\mathcal{U}|\mathcal{V}) \\ R_2 \geq \mathbb{H}(\mathcal{V}|\mathcal{U}) \\ R_1 + R_2 \geq \mathbb{H}(\mathcal{U}\mathcal{V}) \end{array} \right\} \subseteq \mathcal{R}^{\text{SW}}.$$

The converse part of the proof follows from the converse of the source coding theorem and is omitted. \square

Figure 2.6 illustrates the typical shape of the Slepian–Wolf region \mathcal{R}^{SW} . A special case of the Slepian–Wolf problem is when one of the components of the DMS $(\mathcal{U}\mathcal{V}, p_{\mathcal{U}\mathcal{V}})$, say \mathcal{V} , is directly available at the decoder as side information and only \mathcal{U} should be compressed. This problem is known as source coding with side information. The characterization of the minimum compression rate required to reconstruct \mathcal{U} reliably at the decoder follows from Theorem 2.10.

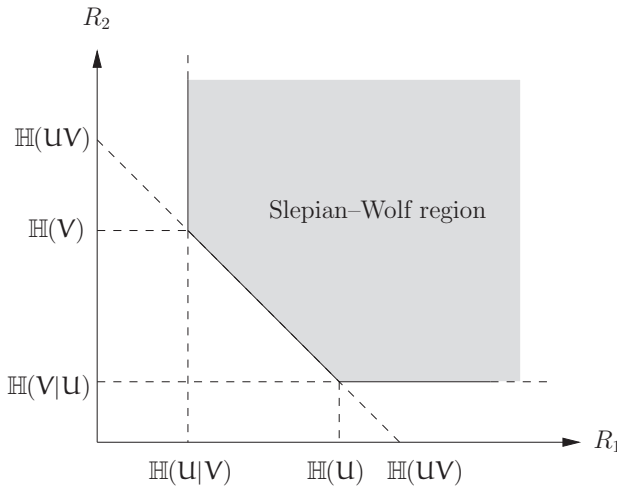


Figure 2.6 The Slepian–Wolf region for a DMS $(\mathcal{UV}, p_{\mathcal{UV}})$.

Corollary 2.4 (Source coding with side information). *Consider a DMS $(\mathcal{UV}, p_{\mathcal{UV}})$ and assume that $(\mathcal{U}, p_{\mathcal{U}})$ should be compressed knowing that $(\mathcal{V}, p_{\mathcal{V}})$ is available as side information at the decoder. Then,*

$$\inf\{R : R \text{ is an achievable compression rate}\} = \mathbb{H}(\mathcal{U}|\mathcal{V}).$$

Corollary 2.4 plays a fundamental role for secret-key agreement in Chapters 4 and 6.

2.3.2 The multiple-access channel

In the previous problem, we assumed that the information generated by multiple sources is transmitted over noiseless channels. If these data are to be communicated over a common noisy channel to a single destination, we call this type of channel a *multiple-access channel* (MAC). As illustrated in Figure 2.7, a discrete memoryless multiple access channel $(\mathcal{X}_1, \mathcal{X}_2, p_{Y|X_1X_2}, \mathcal{Y})$ consists of two finite input alphabets \mathcal{X}_1 and \mathcal{X}_2 , one finite output alphabet \mathcal{Y} , and transition probabilities $p_{Y|X_1X_2}$ such that

$$\forall n \geq 1 \quad \forall (x_1^n, x_2^n, y^n) \in \mathcal{X}_1^n \times \mathcal{X}_2^n \times \mathcal{Y}^n$$

$$p_{Y^n|X_1^nX_2^n}(y^n|x_1^n x_2^n) = \prod_{i=1}^n p_{Y|X_1X_2}(y_i|x_{1,i}, x_{2,i}).$$

Definition 2.23. A $(2^{nR_1}, 2^{nR_2}, n)$ code \mathcal{C}_n for the MAC consists of

- two message sets $\mathcal{M}_1 = \llbracket 1, 2^{nR_1} \rrbracket$ and $\mathcal{M}_2 = \llbracket 1, 2^{nR_2} \rrbracket$;
- two encoding functions, $f_1 : \mathcal{M}_1 \rightarrow \mathcal{X}_1^n$ and $f_2 : \mathcal{M}_2 \rightarrow \mathcal{X}_2^n$, which map a message m_1 or m_2 to a codeword x_1^n or x_2^n ;

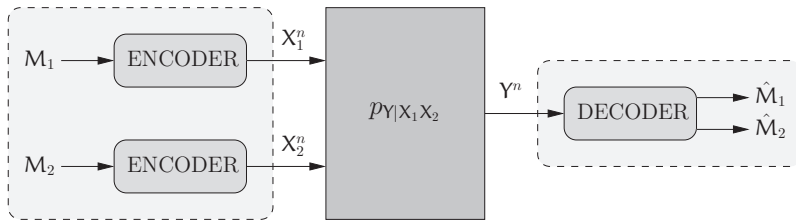


Figure 2.7 Communication over a two-user multiple-access channel.

- a decoding function $g : \mathcal{Y}^n \rightarrow \mathcal{M}_1 \times \mathcal{M}_2 \cup \{?\}$, which maps each channel observation y^n to a message pair $(\hat{m}_1, \hat{m}_2) \in \mathcal{M}_1 \times \mathcal{M}_2$ or an error message ?.

The messages M_1 and M_2 are assumed uniformly distributed in their respective sets, and the performance of a code \mathcal{C}_n is measured in terms of the average probability of error

$$\mathbf{P}_e(\mathcal{C}_n) \triangleq \mathbb{P}[(\hat{M}_1, \hat{M}_2) \neq (M_1, M_2) | \mathcal{C}_n].$$

Definition 2.24. A rate pair (R_1, R_2) is achievable for the MAC if there exists a sequence of $(2^{nR_1}, 2^{nR_2}, n)$ codes $\{\mathcal{C}_n\}_{n \geq 1}$ such that

$$\lim_{n \rightarrow \infty} \mathbf{P}_e(\mathcal{C}_n) = 0.$$

The capacity region of a MAC is defined as

$$\mathcal{C}^{\text{MAC}} \triangleq \text{cl}(\{(R_1, R_2) : (R_1, R_2) \text{ is achievable}\}).$$

The characterization of the capacity region requires the notion of a *convex hull*, which we define below.

Definition 2.25. The convex hull of a set $\mathcal{S} \subseteq \mathbb{R}^n$ is the set

$$\text{co}(\mathcal{S}) \triangleq \left\{ \sum_{i=1}^k \lambda_i x_i^n : k \geq 1, \{\lambda_i\}_k \in [0, 1]^k, \sum_{i=1}^k \lambda_i = 1, \{x_i^n\}_k \in \mathcal{S}^k \right\}.$$

Theorem 2.11 (Ahlswede and Liao). Consider a MAC $(\mathcal{X}_1, \mathcal{X}_2, p_{Y|X_1X_2}, \mathcal{Y})$. For any independent distributions p_{X_1} on \mathcal{X}_1 and p_{X_2} on \mathcal{X}_2 , define the set $\mathcal{R}(p_{X_1}p_{X_2})$ as

$$\mathcal{R}(p_{X_1}p_{X_2}) \triangleq \left\{ (R_1, R_2) : \begin{array}{l} 0 \leq R_1 \leq \mathbb{I}(X_1; Y|X_2) \\ 0 \leq R_2 \leq \mathbb{I}(X_2; Y|X_1) \\ 0 \leq R_1 + R_2 \leq \mathbb{I}(X_1X_2; Y) \end{array} \right\},$$

where the joint distribution of X_1, X_2 , and Y factorizes as $p_{X_1}p_{X_2}p_{Y|X_1X_2}$. Then, the capacity region of a MAC is

$$\mathcal{C}^{\text{MAC}} \triangleq \text{co} \left(\bigcup_{p_{X_1}p_{X_2}} \mathcal{R}(p_{X_1}p_{X_2}) \right).$$

Proof. We provide only the achievability part of the proof, which is similar to the proof of Shannon's channel coding theorem and is based on joint typicality and random

coding. Fix two independent probability distributions, p_{X_1} on \mathcal{X}_1 and p_{X_2} on \mathcal{X}_2 . Let $0 < \epsilon < \mu_{X_1 X_2 Y}$, where

$$\mu_{X_1 X_2 Y} \triangleq \min_{(x_1, x_2, y) \in \mathcal{X}_1 \times \mathcal{X}_2 \times \mathcal{Y}} p_{X_1}(x_1) p_{X_2}(x_2) p_{Y|X_1 X_2}(y|x_1, x_2),$$

and let $n \in \mathbb{N}^*$. Let $R_1 > 0$ and $R_2 > 0$ be rates to be specified later. We construct a $(2^{nR_1}, 2^{nR_2}, n)$ code C_n as follows.

- *Codebook construction.* Construct a codebook for user 1 with $\lceil 2^{nR_1} \rceil$ codewords, labeled $x_1^n(m_1)$ with $m_1 \in \llbracket 1, 2^{nR_1} \rrbracket$, by generating the symbols $x_{1,i}(m_1)$ for $i \in \llbracket 1, n \rrbracket$ and $m_1 \in \llbracket 1, 2^{nR_1} \rrbracket$ independently according to p_{X_1} . Similarly, construct a codebook for user 2 with $\lceil 2^{nR_2} \rceil$ codewords, labeled $x_2^n(m_2)$ with $m_2 \in \llbracket 1, 2^{nR_2} \rrbracket$, by generating the symbols $x_{2,i}(m_2)$ for $i \in \llbracket 1, n \rrbracket$ and $m_2 \in \llbracket 1, 2^{nR_2} \rrbracket$ independently according to p_{X_2} . The codebooks are revealed to all encoders and decoders.
- *Encoder 1.* Given m_1 , transmit $x_1^n(m_1)$.
- *Encoder 2.* Given m_2 , transmit $x_2^n(m_2)$.
- *Decoder.* Given y^n , output (\hat{m}_1, \hat{m}_2) if it is the unique message pair such that $(x_1^n(\hat{m}_1), x_2^n(\hat{m}_2), y^n) \in \mathcal{T}_\epsilon^n(X_1 X_2 Y)$; otherwise, output an error ?.

The random variable that represents the randomly generated code C_n is denoted by C_n and we proceed to bound $\mathbb{E}[\mathbf{P}_e(C_n)]$. By virtue of the symmetry of the random code construction

$$\begin{aligned} \mathbb{E}[\mathbf{P}_e(C_n)] &= \mathbb{E}_{C_n} \left[\mathbb{P} \left[(\hat{M}_1, \hat{M}_2) \neq (M_1, M_2) | C_n \right] \right] \\ &= \mathbb{E}_{C_n} \left[\mathbb{P} \left[(\hat{M}_1, \hat{M}_2) \neq (M_1, M_2) | M_1 = 1, M_2 = 1, C_n \right] \right]. \end{aligned}$$

Therefore, the probability of error can be expressed in terms of the error events

$$\mathcal{E}_{ij} \triangleq \{(X_1^n(i), X_2^n(j), Y^n) \in \mathcal{T}_\epsilon^n(X_1 X_2 Y)\} \quad \text{for } i \in \llbracket 1, 2^{nR_1} \rrbracket \text{ and } j \in \llbracket 1, 2^{nR_2} \rrbracket$$

as

$$\mathbb{E}[\mathbf{P}_e(C_n)] = \mathbb{P} \left[\mathcal{E}_{11}^c \cup \bigcup_{i \neq 1} \mathcal{E}_{i1} \cup \bigcup_{j \neq 1} \mathcal{E}_{1j} \cup \bigcup_{(i,j) \neq (1,1)} \mathcal{E}_{ij} \right].$$

By the union bound,

$$\mathbb{E}[\mathbf{P}_e(C_n)] \leq \mathbb{P}[\mathcal{E}_{11}^c] + \sum_{i \neq 1} \mathbb{P}[\mathcal{E}_{i1}] + \sum_{j \neq 1} \mathbb{P}[\mathcal{E}_{1j}] + \sum_{(i,j) \neq (1,1)} \mathbb{P}[\mathcal{E}_{ij}]. \quad (2.12)$$

By the joint AEP,

$$\mathbb{P}[\mathcal{E}_{11}^c] \leq \delta_\epsilon(n). \quad (2.13)$$

For $i \neq 1$, $X_1^n(i)$ is conditionally independent of Y^n given $X_2^n(1)$; therefore, by Corollary 2.3,

$$\mathbb{P}[\mathcal{E}_{i1}] \leq 2^{-n(\mathbb{I}(X_1; Y|X_2) - \delta(\epsilon))} \quad \text{for } i \neq 1. \quad (2.14)$$

Similarly, we can show

$$\mathbb{P}[\mathcal{E}_{1j}] \leq 2^{-n(\mathbb{I}(X_2; Y|X_1) - \delta(\epsilon))} \quad \text{for } j \neq 1, \quad (2.15)$$

$$\mathbb{P}[\mathcal{E}_{ij}] \leq 2^{-n(\mathbb{I}(X_1 X_2; Y) - \delta(\epsilon))} \quad \text{for } i \neq 1 \text{ and } j \neq 1. \quad (2.16)$$

On substituting (2.13)–(2.16) into (2.12), we obtain

$$\begin{aligned} \mathbb{E}[\mathbf{P}_e(C_n)] &\leq \delta_\epsilon(n) + \lceil 2^{nR_1} \rceil 2^{-n(\mathbb{I}(X_1; Y|X_2) - \delta(\epsilon))} + \lceil 2^{nR_2} \rceil 2^{-n(\mathbb{I}(X_2; Y|X_1) - \delta(\epsilon))} \\ &\quad + \lceil 2^{nR_1} \rceil \lceil 2^{nR_2} \rceil 2^{-n(\mathbb{I}(X_1 X_2; Y) - \delta(\epsilon))}. \end{aligned}$$

Hence, if we choose R_1 and R_2 to satisfy

$$R_1 < \mathbb{I}(X_1; Y|X_2) - \delta(\epsilon),$$

$$R_2 < \mathbb{I}(X_2; Y|X_1) - \delta(\epsilon),$$

$$R_1 + R_2 < \mathbb{I}(X_1 X_2; Y) - \delta(\epsilon),$$

we obtain $\mathbb{E}[\mathbf{P}_e(C_n)] \leq \delta_\epsilon(n)$. By applying the selection lemma to the random variable C_n and the function \mathbf{P}_e , we conclude that there exists a $(2^{nR_1}, 2^{nR_2}, n)$ code C_n such that $\mathbf{P}_e(C_n) \leq \delta_\epsilon(n)$. Since ϵ can be chosen arbitrarily small and since the distributions p_{X_1} and p_{X_2} are arbitrary, we conclude that

$$\bigcup_{p_{X_1} p_{X_2}} \left\{ (R_1, R_2): \begin{array}{l} 0 \leq R_1 \leq \mathbb{I}(X_1; Y|X_2) \\ 0 \leq R_2 \leq \mathbb{I}(X_2; Y|X_1) \\ 0 \leq R_1 + R_2 \leq \mathbb{I}(X_1 X_2; Y) \end{array} \right\} \subseteq \mathcal{C}^{\text{MAC}}$$

is achievable. Finally, it can be shown that time-sharing between different codes achieves the entire convex hull [3, Section 15.3]. We refer the reader to [3, 6] for the converse part of the proof. \square

The typical shape of the region $\mathcal{R}(p_{X_1} p_{X_2})$ is illustrated in Figure 2.8. The boundaries of the capacity region can be explained in a very intuitive way. When encoder 1 views the signals sent by encoder 2 as noise, its maximum achievable rate is on the order of $\mathbb{I}(X_1; Y)$, which is a direct consequence of the channel coding theorem. Then, the decoder can estimate the codeword x_1^n and subtract it from the channel output sequence y_1^n , thus allowing encoder 2 to achieve a maximum rate on the order of $\mathbb{I}(X_2; Y|X_1)$. This procedure is sometimes called *successive cancellation* and leads to the upper corner point of the region. The lower corner point corresponds to the symmetric case, in which encoder 2 views the signals sent by encoder 1 as noise.

2.3.3 The broadcast channel

While a multiple-access channel considers multiple sources and one destination, the *broadcast channel* (BC for short) considers a single information source that transmits to multiple users. Applications of the BC include the downlink channel of a satellite or of a base station in a mobile communication network, and the wiretap channel which is studied in detail in Chapter 3 and Chapter 5. As illustrated in Figure 2.9, a discrete

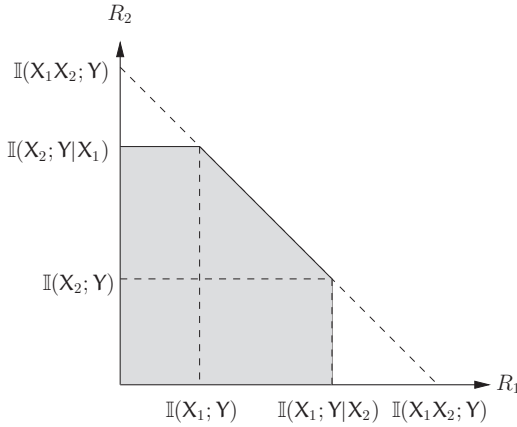


Figure 2.8 Typical shape of the rate region $\mathcal{R}(p_{X_1} p_{X_2})$ of the multiple-access channel for fixed input distributions p_{X_1} and p_{X_2} .

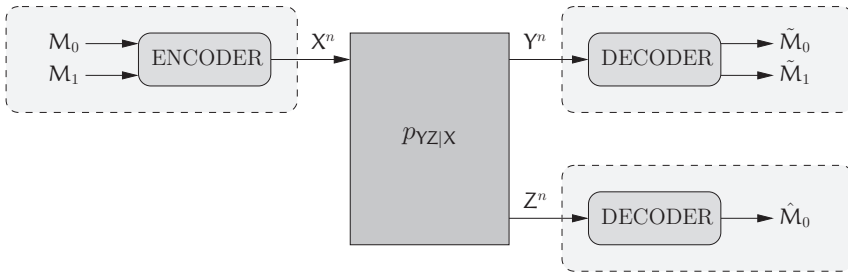


Figure 2.9 Communication over a two-user broadcast channel.

memoryless two-user broadcast channel $(\mathcal{X}, p_{Y|Z|X}, \mathcal{Y}, \mathcal{Z})$ consists of a finite input alphabet \mathcal{X} , two finite output alphabets \mathcal{Y} and \mathcal{Z} , and transition probabilities $p_{Y|Z|X}$ such that

$$\forall n \geq 1 \quad \forall (x^n, y^n, z^n) \in \mathcal{X}^n \times \mathcal{Y}^n \times \mathcal{Z}^n$$

$$p_{Y^n Z^n | X^n}(y^n, z^n | x^n) = \prod_{i=1}^n p_{Y|Z|X}(y_i, z_i | x_i).$$

We assume that the transmitter wants to send a *common message* M_0 to both receivers and a *private message* M_1 to the receiver observing Y^n . The receiver observing Z^n is called a “weak” user, while the receiver observing Y^n is called the “strong” user.

Definition 2.26. A $(2^{nR_0}, 2^{nR_1}, n)$ code C_n for the BC consists of

- two message sets $\mathcal{M}_0 = [1, 2^{nR_0}]$ and $\mathcal{M}_1 = [1, 2^{nR_1}]$;
- an encoding function $f : \mathcal{M}_0 \times \mathcal{M}_1 \rightarrow \mathcal{X}^n$, which maps a message pair (m_0, m_1) to a codeword x^n ;

- a decoding function $g : \mathcal{Y}^n \rightarrow (\mathcal{M}_0 \times \mathcal{M}_1) \cup \{?\}$, which maps each channel observation y^n to a message pair $(\tilde{m}_0, \tilde{m}_1) \in \mathcal{M}_0 \times \mathcal{M}_1$ or an error message ?;
- a decoding function $h : \mathcal{Z}^n \rightarrow \mathcal{M}_0 \cup \{?\}$, which maps each channel observation z^n to a message $\hat{m}_0 \in \mathcal{M}_0$ or an error message ?.

Messages M_0 and M_1 are assumed uniformly distributed in their respective sets and the performance of a code C_n is measured in terms of the average probability of error

$$\mathbf{P}_e(C_n) \triangleq \mathbb{P} \left[\hat{M}_0 \neq M_0 \text{ or } (\tilde{M}_0, \tilde{M}_1) \neq (M_0, M_1) \mid C_n \right].$$

Definition 2.27. A rate pair (R_0, R_1) is achievable for the BC if there exists a sequence of $(2^{nR_0}, 2^{nR_1}, n)$ codes $\{C_n\}_{n \geq 1}$ such that

$$\lim_{n \rightarrow \infty} \mathbf{P}_e(C_n) = 0.$$

The capacity region of a BC is defined as

$$\mathcal{C}^{\text{BC}} \triangleq \text{cl}(\{(R_0, R_1) : (R_0, R_1) \text{ is achievable}\}).$$

As in many other fundamental problems of network information theory, determining the capacity region of the broadcast channel turns out to be a very difficult task. Therefore, we provide only an achievable rate region, which, in general, is strictly smaller than the capacity region.

Theorem 2.12 (Bergsman and Gallager). Consider a BC $(\mathcal{X}, p_{Y|Z|X}, \mathcal{Y}, \mathcal{Z})$. For any joint distribution $p_{\mathcal{U}X}$ on $\mathcal{U} \times \mathcal{X}$, define the set $\mathcal{R}(p_{\mathcal{U}X})$ as

$$\mathcal{R}(p_{\mathcal{U}X}) \triangleq \left\{ (R_0, R_1) : \begin{array}{l} 0 \leq R_0 \leq \min(\mathbb{I}(\mathcal{U}; \mathcal{Y}), \mathbb{I}(\mathcal{U}; \mathcal{Z})) \\ 0 \leq R_1 \leq \mathbb{I}(\mathcal{X}; \mathcal{Y}|\mathcal{U}) \end{array} \right\},$$

where the joint distribution of $\mathcal{U}, \mathcal{X}, \mathcal{Y}$, and \mathcal{Z} factorizes as $p_{\mathcal{U}X} p_{Y|Z|X}$. Then,

$$\mathcal{R}^{\text{BC}} \triangleq \text{co} \left(\bigcup_{p_{\mathcal{U}X}} \mathcal{R}(p_{\mathcal{U}X}) \right) \subseteq \mathcal{C}^{\text{BC}}.$$

In addition, the cardinality of the auxiliary random variable \mathcal{U} can be limited to $|\mathcal{U}| \leq \min(|\mathcal{X}|, |\mathcal{Y}|, |\mathcal{Z}|)$.

Proof. The proof that $\mathcal{R}^{\text{BC}} \subseteq \mathcal{C}^{\text{BC}}$ is based on joint typicality, random coding, and a code construction called *superposition coding*. As illustrated in Figure 2.10, the idea of superposition coding is to create a codebook with $\lceil 2^{nR_0} \rceil$ codewords for the weakest user and to superpose a codebook with $\lceil 2^{nR_1} \rceil$ codewords for the strongest user to every codeword. The codewords u^n are often called “cloud centers,” while the codewords x^n are called “satellite codewords.” Formally, fix a joint probability distribution $p_{\mathcal{U}X}$ on $\mathcal{U} \times \mathcal{X}$. Let $0 < \epsilon < \mu_{XYU}$, where $\mu_{XYU} \triangleq \min p_{XU}(x, u) p_{Y|XU}(y|x, u)$ and let $n \in \mathbb{N}^*$. Let $R_0 > 0$ and $R_1 > 0$ be rates to be specified later. We construct a $(2^{nR_0}, 2^{nR_1}, n)$ code C_n as follows.

- *Codebook construction.* Construct a codebook with $\lceil 2^{nR_0} \rceil$ codewords, labeled $u^n(m_0)$ with $m_0 \in \llbracket 1, 2^{nR_0} \rrbracket$, by generating the symbols $u_i(m_0)$ for $i \in \llbracket 1, n \rrbracket$ and

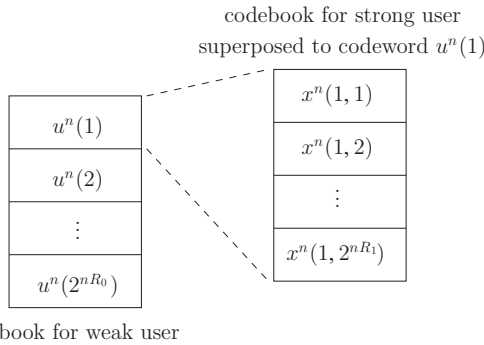


Figure 2.10 Superposition coding. A codebook for the strong user is superposed to every codeword of the codebook for the weak user.

$m_0 \in \llbracket 1, 2^{nR_0} \rrbracket$ independently according to p_U . For each $u^n(m_0)$ with $m_0 \in \llbracket 1, 2^{nR_0} \rrbracket$, generate another codebook with $\lceil 2^{nR_1} \rceil$ codewords, labeled $x^n(m_0, m_1)$ with $m_1 \in \llbracket 1, 2^{nR_1} \rrbracket$, by generating the symbols $x_i(m_0, m_1)$ for $i \in \llbracket 1, n \rrbracket$ and $m_1 \in \llbracket 1, 2^{nR_1} \rrbracket$ independently according to $p_{X|U=u_i(m_0)}$. The codebooks are revealed to the encoder and both decoders.

- *Encoder*: Given (m_0, m_1) , transmit $x^n(m_0, m_1)$.
- *Decoder for weak user*: Given z^n , output \hat{m}_0 if it is the unique message such that $(u^n(\hat{m}_0), z^n) \in \mathcal{T}_\epsilon^n(\mathcal{U}Z)$; otherwise, output an error ?.
- *Decoder for strong user*: Given y^n , output $(\tilde{m}_0, \tilde{m}_1)$ if it is the unique message pair such that $(u^n(\tilde{m}_0), y^n) \in \mathcal{T}_\epsilon^n(\mathcal{U}Y)$ and $(u^n(\tilde{m}_0), x^n(\tilde{m}_0, \tilde{m}_1), y^n) \in \mathcal{T}_\epsilon^n(\mathcal{U}XY)$; otherwise, output an error ?.

The random variable that denotes the randomly generated codebook \mathcal{C}_n is denoted by C_n , and we proceed to bound $\mathbb{E}[\mathbf{P}_e(C_n)]$. From the symmetry of the random code construction, notice that

$$\begin{aligned} \mathbb{E}[\mathbf{P}_e(C_n)] &= \mathbb{E}_{C_n} \left[\mathbb{P} \left[\hat{M}_0 \neq M_0 \text{ or } (\tilde{M}_0, \tilde{M}_1) \neq (M_0, M_1) \mid C_n \right] \right] \\ &= \mathbb{E}_{C_n} \left[\mathbb{P} \left[\hat{M}_0 \neq M_0 \text{ or } (\tilde{M}_0, \tilde{M}_1) \neq (M_0, M_1) \mid M_0 = 1, M_1 = 1, C_n \right] \right]. \end{aligned}$$

Therefore, $\mathbb{E}[\mathbf{P}_e(C_n)]$ can be expressed in terms of the events

$$\begin{aligned} \mathcal{E}_i &= \{(\mathcal{U}^n(i), Y^n) \in \mathcal{T}_\epsilon^n(\mathcal{U}Y)\} \quad \text{for } i \in \llbracket 1, 2^{nR_0} \rrbracket, \\ \mathcal{F}_i &= \{(\mathcal{U}^n(i), Z^n) \in \mathcal{T}_\epsilon^n(\mathcal{U}Z)\} \quad \text{for } i \in \llbracket 1, 2^{nR_0} \rrbracket, \\ \mathcal{G}_{ij} &= \{(\mathcal{U}^n(i), X^n(j), Y^n) \in \mathcal{T}_\epsilon^n(\mathcal{U}XY)\} \quad \text{for } i \in \llbracket 1, 2^{nR_0} \rrbracket \text{ and } j \in \llbracket 1, 2^{nR_1} \rrbracket \end{aligned}$$

as

$$\mathbb{E}[\mathbf{P}_e(C_n)] = \mathbb{P} \left[\mathcal{E}_1^c \cup \bigcup_{i \neq 1} \mathcal{E}_i \cup \mathcal{F}_1^c \cup \bigcup_{i \neq 1} \mathcal{F}_i \cup \bigcup_{j \neq 1} \mathcal{G}_{1j} \right]$$

and, by the union bound,

$$\mathbb{E}[\mathbf{P}_e(C_n)] \leq \mathbb{P}[\mathcal{E}_1^c] + \sum_{i \neq 1} \mathbb{P}[\mathcal{E}_i] + \mathbb{P}[\mathcal{F}_1^c] + \sum_{i \neq 1} \mathbb{P}[\mathcal{F}_i] + \sum_{j \neq 1} \mathbb{P}[\mathcal{G}_{1j}]. \quad (2.17)$$

By the joint AEP,

$$\mathbb{P}[\mathcal{E}_1^c] \leq \delta_\epsilon(n) \quad \text{and} \quad \mathbb{P}[\mathcal{F}_1^c] \leq \delta_\epsilon(n). \quad (2.18)$$

For $i \neq 1$, notice that $\mathbf{U}^n(i)$ is independent of \mathbf{Y}^n and \mathbf{Z}^n ; therefore, by Corollary 2.2,

$$\mathbb{P}[\mathcal{E}_i] \leq 2^{-n(\mathbb{I}(\mathbf{U}; \mathbf{Y}) - \delta(\epsilon))} \quad \text{and} \quad \mathbb{P}[\mathcal{F}_i] \leq 2^{-n(\mathbb{I}(\mathbf{U}; \mathbf{Z}) - \delta(\epsilon))} \quad \text{for } i \neq 1. \quad (2.19)$$

For $j \neq 1$, $\mathbf{X}^n(1, j)$ is conditionally independent of \mathbf{Z}^n given $\mathbf{U}^n(1)$; therefore, by Corollary 2.3,

$$\mathbb{P}[\mathcal{G}_{1j}] \leq 2^{-n(\mathbb{I}(\mathbf{X}; \mathbf{Y}|\mathbf{U}) - \delta(\epsilon))} \quad \text{for } j \neq 1. \quad (2.20)$$

On substituting (2.18), (2.19), and (2.20) into (2.17), we obtain

$$\begin{aligned} \mathbb{E}[\mathbf{P}_e(C_n)] &\leq \delta_\epsilon(n) + \lceil 2^{nR_0} \rceil 2^{-n(\mathbb{I}(\mathbf{U}; \mathbf{Y}) - \delta(\epsilon))} + \lceil 2^{nR_0} \rceil 2^{-n(\mathbb{I}(\mathbf{U}; \mathbf{Z}) - \delta(\epsilon))} \\ &\quad + \lceil 2^{nR_1} \rceil 2^{-n(\mathbb{I}(\mathbf{X}; \mathbf{Y}|\mathbf{U}) - \delta(\epsilon))}. \end{aligned}$$

Hence, if we choose the rates R_0 and R_1 to satisfy

$$R_0 < \min(\mathbb{I}(\mathbf{U}; \mathbf{Y}), \mathbb{I}(\mathbf{U}; \mathbf{Z})) - \delta(\epsilon),$$

$$R_1 < \mathbb{I}(\mathbf{X}; \mathbf{Y}|\mathbf{U}) - \delta(\epsilon),$$

we obtain $\mathbb{E}[C_n] \leq \delta_\epsilon(n)$. By applying the selection lemma to the random variable C_n and the function \mathbf{P}_e , we conclude that there exists a $(2^{nR_0}, 2^{nR_1}, n)$ code C_n such that $\mathbf{P}_e(C_n) \leq \delta_\epsilon(n)$. Since ϵ can be chosen arbitrarily small and since the distribution $p_{\mathbf{U}\mathbf{X}}$ is arbitrary, we conclude that

$$\left\{ (R_0, R_1): \begin{array}{l} 0 \leq R_0 \leq \min(\mathbb{I}(\mathbf{U}; \mathbf{Y}), \mathbb{I}(\mathbf{U}; \mathbf{Z})) \\ 0 \leq R_1 \leq \mathbb{I}(\mathbf{X}; \mathbf{Y}|\mathbf{U}) \end{array} \right\} \subseteq \mathcal{C}^{\text{BC}}.$$

Since $p_{\mathbf{U}\mathbf{X}}$ is arbitrary and since it is possible to perform time-sharing, the theorem follows. The bound for the cardinality of the random variable \mathbf{U} follows from Caratheodory's theorem, and we refer the reader to [3] for details. \square

2.4 Bibliographical notes

The definitions of typical sequences and their properties are those described in the textbooks [3, 4, 6]. The notion of d-separation is a known result in statistical inference, and we have used the definition provided by Kramer [8].

There are several ways of proving the channel coding theorem, see for instance [2, 3, 4, 5], and our presentation is based on the approach in [3, 4]. The Slepian–Wolf theorem was established by Slepian and Wolf in [9]. Additional examples of results proved with random binning with side information can be found in [10] and [11].

The capacity of the two-user multiple-access channel with independent messages was obtained independently by Ahlswede [12] and Liao [13]. The broadcast channel model was proposed by Cover [14]. Bergmans [15] and Gallager [16] proved that the achievable region obtained in this chapter is the capacity region of a subclass of broadcast channels called *physically degraded* broadcast channels. Surveys of known results about network information theory can be found in Cover's survey paper [17] or Kramer's monograph [6].

