

## Chapter 12

---

# Information Theory and Statistics

---

We now explore the relationship between information theory and statistics. We begin by describing the method of types, which is a powerful technique in large deviation theory. We use the method of types to calculate the probability of rare events and to show the existence of universal source codes. We also consider the problem of testing hypotheses and derive the best possible error exponents for such tests (Stein's lemma). Finally, we treat the estimation of the parameters of a distribution and describe the role of Fisher information.

### 12.1 THE METHOD OF TYPES

The AEP for discrete random variables (Chapter 3) focuses our attention on a small subset of typical sequences. The method of types is an even more powerful procedure in which we consider the sequences that have the same empirical distribution. With this restriction, we can derive strong bounds on the number of sequences of a particular empirical distribution and the probability of each sequence in this set. It is then possible to derive strong error bounds for the channel coding theorem and prove a variety of rate-distortion results. The method of types was fully developed by Csiszár and Körner [83], who obtained most of their results from this point of view.

Let  $X_1, X_2, \dots, X_n$  be a sequence of  $n$  symbols from an alphabet  $\mathcal{X} = \{a_1, a_2, \dots, a_{|\mathcal{X}|}\}$ . We will use the notation  $x^n$  and  $\mathbf{x}$  interchangeably to denote a sequence  $x_1, x_2, \dots, x_n$ .

**Definition:** The *type*  $P_{\mathbf{x}}$  (or empirical probability distribution) of a sequence  $x_1, x_2, \dots, x_n$  is the relative proportion of occurrences of each

symbol of  $\mathcal{X}$ , i.e.,  $P_{\mathbf{x}}(a) = N(a|\mathbf{x})/n$  for all  $a \in \mathcal{X}$ , where  $N(a|\mathbf{x})$  is the number of times the symbol  $a$  occurs in the sequence  $\mathbf{x} \in \mathcal{X}^n$ .

The type of a sequence  $\mathbf{x}$  is denoted as  $P_{\mathbf{x}}$ . It is a probability mass function on  $\mathcal{X}$ . (Note that in this chapter, we will use capital letters to denote types and distributions. We will also loosely use the word "distribution" to mean a probability mass function.)

**Definition:** Let  $\mathcal{P}_n$  denote the set of types with denominator  $n$ .

For example, if  $\mathcal{X} = \{0, 1\}$ , then the set of possible types with denominator  $n$  is

$$\mathcal{P}_n = \left\{ (P(0), P(1)) : \left( \frac{0}{n}, \frac{n}{n} \right), \left( \frac{1}{n}, \frac{n-1}{n} \right), \dots, \left( \frac{n}{n}, \frac{0}{n} \right) \right\}. \quad (12.1)$$

**Definition:** If  $P \in \mathcal{P}_n$ , then the set of sequences of length  $n$  and type  $P$  is called the *type class* of  $P$ , denoted  $T(P)$ , i.e.,

$$T(P) = \{ \mathbf{x} \in \mathcal{X}^n : P_{\mathbf{x}} = P \}. \quad (12.2)$$

The type class is sometimes called the composition class of  $P$ .

**Example 12.1.1:** Let  $\mathcal{X} = \{1, 2, 3\}$ , a ternary alphabet. Let  $\mathbf{x} = 11321$ . Then the type  $P_{\mathbf{x}}$  is

$$P_{\mathbf{x}}(1) = \frac{3}{5}, \quad P_{\mathbf{x}}(2) = \frac{1}{5}, \quad P_{\mathbf{x}}(3) = \frac{1}{5}. \quad (12.3)$$

The type class of  $P_{\mathbf{x}}$  is the set of all sequences of length 5 with three 1's, one 2 and one 3. There are 20 such sequences, and

$$T(P_{\mathbf{x}}) = \{11123, 11132, 11213, \dots, 32111\}. \quad (12.4)$$

The number of elements in  $T(P)$  is

$$|T(P)| = \binom{5}{3, 1, 1} = \frac{5!}{3!1!1!} = 20. \quad (12.5)$$

The essential power of the method of types arises from the following theorem, which shows that the number of types is at most polynomial in  $n$ .

**Theorem 12.1.1:**

$$|\mathcal{P}_n| \leq (n+1)^{|\mathcal{X}|}. \quad (12.6)$$

**Proof:** There are  $|\mathcal{X}|$  components in the vector that specifies  $P_{\mathbf{x}}$ . The numerator in each component can take on only  $n + 1$  values. So there are at most  $(n + 1)^{|\mathcal{X}|}$  choices for the type vector. Of course, these choices are not independent (for example, the last choice is fixed by the others). But this is a sufficiently good upper bound for our needs.  $\square$

The crucial point here is that there are only a polynomial number of types of length  $n$ . Since the number of sequences is exponential in  $n$ , it follows that at least one type has exponentially many sequences in its type class. In fact, the largest type class has essentially the same number of elements as the entire set of sequences, to first order in the exponent.

Now, we will assume that the sequence  $X_1, X_2, \dots, X_n$  is drawn i.i.d. according to a distribution  $Q(x)$ . All sequences with the same type will have the same probability, as shown in the following theorem. Let  $Q^n(x^n) = \prod_{i=1}^n Q(x_i)$  denote the product distribution associated with  $Q$ .

**Theorem 12.1.2:** *If  $X_1, X_2, \dots, X_n$  are drawn i.i.d. according to  $Q(x)$ , then the probability of  $\mathbf{x}$  depends only on its type and is given by*

$$Q^n(\mathbf{x}) = 2^{-n(H(P_{\mathbf{x}}) + D(P_{\mathbf{x}} \| Q))}. \quad (12.7)$$

**Proof:**

$$Q^n(\mathbf{x}) = \prod_{i=1}^n Q(x_i) \quad (12.8)$$

$$= \prod_{a \in \mathcal{X}} Q(a)^{N(a|\mathbf{x})} \quad (12.9)$$

$$= \prod_{a \in \mathcal{X}} Q(a)^{nP_{\mathbf{x}}(a)} \quad (12.10)$$

$$= \prod_{a \in \mathcal{X}} 2^{nP_{\mathbf{x}}(a) \log Q(a)} \quad (12.11)$$

$$= \prod_{a \in \mathcal{X}} 2^{n(P_{\mathbf{x}}(a) \log Q(a) - P_{\mathbf{x}}(a) \log P_{\mathbf{x}}(a) + P_{\mathbf{x}}(a) \log P_{\mathbf{x}}(a))} \quad (12.12)$$

$$= 2^{n \sum_{a \in \mathcal{X}} (-P_{\mathbf{x}}(a) \log \frac{P_{\mathbf{x}}(a)}{Q(a)} + P_{\mathbf{x}}(a) \log P_{\mathbf{x}}(a))} \quad (12.13)$$

$$= 2^{n(-D(P_{\mathbf{x}} \| Q) - H(P_{\mathbf{x}}))}. \quad \square \quad (12.14)$$

**Corollary:** *If  $\mathbf{x}$  is in the type class of  $Q$ , then*

$$Q^n(\mathbf{x}) = 2^{-nH(Q)}. \quad (12.15)$$

**Proof:** If  $\mathbf{x} \in T(Q)$ , then  $P_{\mathbf{x}} = Q$ , which can be substituted into (12.14).  $\square$

**Example 12.1.2:** The probability that a fair die produces a particular sequence of length  $n$  with precisely  $n/6$  occurrences of each face ( $n$  is a multiple of 6) is  $2^{-nH(\frac{1}{6}, \frac{1}{6}, \dots, \frac{1}{6})} = 6^{-n}$ . This is obvious. However, if the die has a probability mass function  $(\frac{1}{3}, \frac{1}{3}, \frac{1}{6}, \frac{1}{12}, \frac{1}{12}, 0)$ , the probability of observing a particular sequence with precisely these frequencies is precisely  $2^{-nH(\frac{1}{3}, \frac{1}{3}, \frac{1}{6}, \frac{1}{12}, \frac{1}{12}, 0)}$  for  $n$  a multiple of 12. This is more interesting.

We now give an estimate of the size of a type class  $T(P)$ .

**Theorem 12.1.3** (*Size of a type class  $T(P)$* ): For any type  $P \in \mathcal{P}_n$ ,

$$\frac{1}{(n+1)^{|\mathcal{X}|}} 2^{nH(P)} \leq |T(P)| \leq 2^{nH(P)}. \quad (12.16)$$

**Proof:** The exact size of  $T(P)$  is easy to calculate. It is a simple combinatorial problem—the number of ways of arranging  $nP(a_1)$ ,  $nP(a_2), \dots, nP(a_{|\mathcal{X}|})$  objects in a sequence, which is

$$|T(P)| = \binom{n}{nP(a_1), nP(a_2), \dots, nP(a_{|\mathcal{X}|})}. \quad (12.17)$$

This value is hard to manipulate, so we derive simple exponential bounds on its value.

We suggest two alternative proofs for the exponential bounds.

The first proof uses Stirling's formula [110] to bound the factorial function, and after some algebra, we can obtain the bounds of the theorem.

We give an alternative proof. We first prove the upper bound. Since a type class must have probability  $\leq 1$ , we have

$$1 \geq P^n(T(P)) \quad (12.18)$$

$$= \sum_{\mathbf{x} \in T(P)} P^n(\mathbf{x}) \quad (12.19)$$

$$= \sum_{\mathbf{x} \in T(P)} 2^{-nH(P)} \quad (12.20)$$

$$= |T(P)| 2^{-nH(P)}, \quad (12.21)$$

using Theorem 12.1.2. Thus

$$|T(P)| \leq 2^{nH(P)}. \quad (12.22)$$

Now for the lower bound. We first prove that the type class  $T(P)$  has the highest probability among all type classes under the probability distribution  $P$ , i.e.,

$$P^n(T(P)) \geq P^n(T(\hat{P})), \quad \text{for all } \hat{P} \in \mathcal{P}_n. \quad (12.23)$$

We lower bound the ratio of probabilities,

$$\frac{P^n(T(P))}{P^n(T(\hat{P}))} = \frac{|T(P)| \prod_{a \in \mathcal{X}} P(a)^{nP(a)}}{|T(\hat{P})| \prod_{a \in \mathcal{X}} P(a)^{n\hat{P}(a)}} \quad (12.24)$$

$$= \frac{(\binom{n}{nP(a_1), nP(a_2), \dots, nP(a_{|\mathcal{X}|})}) \prod_{a \in \mathcal{X}} P(a)^{nP(a)}}{(\binom{n}{n\hat{P}(a_1), n\hat{P}(a_2), \dots, n\hat{P}(a_{|\mathcal{X}|})}) \prod_{a \in \mathcal{X}} P(a)^{n\hat{P}(a)}} \quad (12.25)$$

$$= \prod_{a \in \mathcal{X}} \frac{(n\hat{P}(a))!}{(nP(a))!} P(a)^{n(P(a) - \hat{P}(a))}. \quad (12.26)$$

Now using the simple bound (easy to prove by separately considering the cases  $m \geq n$  and  $m < n$ )

$$\frac{m!}{n!} \geq n^{m-n}, \quad (12.27)$$

we obtain

$$\frac{P^n(T(P))}{P^n(T(\hat{P}))} \geq \prod_{a \in \mathcal{X}} (nP(a))^{n\hat{P}(a) - nP(a)} P(a)^{n(P(a) - \hat{P}(a))} \quad (12.28)$$

$$= \prod_{a \in \mathcal{X}} n^{n(\hat{P}(a) - P(a))} \quad (12.29)$$

$$= n^{n(\sum_{a \in \mathcal{X}} \hat{P}(a) - \sum_{a \in \mathcal{X}} P(a))} \quad (12.30)$$

$$= n^{n(1-1)} \quad (12.31)$$

$$= 1. \quad (12.32)$$

Hence  $P^n(T(P)) \geq P^n(T(\hat{P}))$ . The lower bound now follows easily from this result, since

$$1 = \sum_{Q \in \mathcal{P}_n} P^n(T(Q)) \quad (12.33)$$

$$\leq \sum_{Q \in \mathcal{P}_n} \max_Q P^n(T(Q)) \quad (12.34)$$

$$= \sum_{Q \in \mathcal{P}_n} P^n(T(P)) \quad (12.35)$$

$$\leq (n+1)^{|\mathcal{X}|} P^n(T(P)) \quad (12.36)$$

$$= (n+1)^{|\mathcal{X}|} \sum_{\mathbf{x} \in T(P)} P^n(\mathbf{x}) \quad (12.37)$$

$$= (n+1)^{|\mathcal{X}|} \sum_{\mathbf{x} \in T(P)} 2^{-nH(P)} \quad (12.38)$$

$$= (n+1)^{|\mathcal{X}|} |T(P)| 2^{-nH(P)}, \quad (12.39)$$

where (12.36) follows from Theorem 12.1.1 and (12.38) follows from Theorem 12.1.2.  $\square$

We give a slightly better approximation for the binary case.

**Example 12.1.3 (Binary alphabet):** In this case, the type is defined by the number of 1's in the sequence, and the size of the type class is therefore  $\binom{n}{k}$ . We show that

$$\frac{1}{n+1} 2^{nH(\frac{k}{n})} \leq \binom{n}{k} \leq 2^{nH(\frac{k}{n})}. \quad (12.40)$$

These bounds can be proved using Stirling's approximation for the factorial function. But we provide a more intuitive proof below.

We first prove the upper bound. From the binomial formula, for any  $p$ ,

$$\sum_{k=0}^n \binom{n}{k} p^k (1-p)^{n-k} = 1. \quad (12.41)$$

Since all the terms of the sum are positive for  $0 \leq p \leq 1$ , each of the terms is less than 1. Setting  $p = \frac{k}{n}$  and taking the  $k$ th term, we get

$$1 \geq \binom{n}{k} \left(\frac{k}{n}\right)^k \left(1 - \frac{k}{n}\right)^{n-k} \quad (12.42)$$

$$= \binom{n}{k} 2^{k \log \frac{k}{n} + (n-k) \log \frac{n-k}{n}} \quad (12.43)$$

$$= \binom{n}{k} 2^{n(\frac{k}{n} \log \frac{k}{n} + \frac{n-k}{n} \log \frac{n-k}{n})} \quad (12.44)$$

$$= \binom{n}{k} 2^{-nH(\frac{k}{n})}. \quad (12.45)$$

Hence

$$\binom{n}{k} \leq 2^{nH(\frac{k}{n})}. \quad (12.46)$$

For the lower bound, let  $S$  be a random variable with a binomial distribution with parameters  $n$  and  $p$ . The most likely value of  $S$  is  $S = \langle np \rangle$ . This can be easily verified from the fact that

$$\frac{P(S = i + 1)}{P(S = i)} = \frac{n - i}{i + 1} \frac{p}{1 - p} \quad (12.47)$$

and considering the cases when  $i < np$  and when  $i > np$ . Then, since there are  $n + 1$  terms in the binomial sum,

$$1 = \sum_{k=0}^n \binom{n}{k} p^k (1-p)^{n-k} \leq (n+1) \max_k \binom{n}{k} p^k (1-p)^{n-k} \quad (12.48)$$

$$= (n+1) \binom{n}{\langle np \rangle} p^{\langle np \rangle} (1-p)^{n-\langle np \rangle}. \quad (12.49)$$

Now let  $p = \frac{k}{n}$ . Then we have

$$1 \leq (n+1) \binom{n}{k} \left(\frac{k}{n}\right)^k \left(1 - \frac{k}{n}\right)^{n-k}, \quad (12.50)$$

which by the arguments in (12.45) is equivalent to

$$\frac{1}{n+1} \leq \binom{n}{k} 2^{-nH(\frac{k}{n})}, \quad (12.51)$$

or

$$\binom{n}{k} \geq \frac{2^{nH(\frac{k}{n})}}{n+1}. \quad (12.52)$$

Combining the two results, we see that

$$\binom{n}{k} \doteq 2^{nH(\frac{k}{n})}. \quad (12.53)$$

**Theorem 12.1.4** (*Probability of type class*): For any  $P \in \mathcal{P}_n$  and any distribution  $Q$ , the probability of the type class  $T(P)$  under  $Q^n$  is  $2^{-nD(P\|Q)}$  to first order in the exponent. More precisely,

$$\frac{1}{(n+1)^{|\mathcal{X}|}} 2^{-nD(P\|Q)} \leq Q^n(T(P)) \leq 2^{-nD(P\|Q)}. \quad (12.54)$$

**Proof:** We have

$$Q^n(T(P)) = \sum_{\mathbf{x} \in T(P)} Q^n(\mathbf{x}) \quad (12.55)$$

$$= \sum_{\mathbf{x} \in T(P)} 2^{-n(D(P\|Q) + H(P))} \quad (12.56)$$

$$= |T(P)| 2^{-n(D(P\|Q) + H(P))}, \quad (12.57)$$

by Theorem 12.1.2. Using the bounds on  $|T(P)|$  derived in Theorem 12.1.3, we have

$$\frac{1}{(n+1)^{|\mathcal{X}|}} 2^{-nD(P\|Q)} \leq Q^n(T(P)) \leq 2^{-nD(P\|Q)}. \quad \square \quad (12.58)$$

We can summarize the basic theorems concerning types in four equations:

$$|\mathcal{P}_n| \leq (n+1)^{|\mathcal{X}|}, \quad (12.59)$$

$$Q^n(\mathbf{x}) = 2^{-n(D(P_{\mathbf{x}}\|Q) + H(P_{\mathbf{x}}))}, \quad (12.60)$$

$$|T(P)| \doteq 2^{nH(P)}, \quad (12.61)$$

$$Q^n(T(P)) \doteq 2^{-nD(P\|Q)}. \quad (12.62)$$

These equations state that there are only a polynomial number of types and that there are an exponential number of sequences of each type. We also have an exact formula for the probability of any sequence of type  $P$  under distribution  $Q$  and an approximate formula for the probability of a type class.

These equations allow us to calculate the behavior of long sequences based on the properties of the type of the sequence. For example, for long sequences drawn i.i.d. according to some distribution, the type of the sequence is close to the distribution generating the sequence, and we can use the properties of this distribution to estimate the properties of the sequence. Some of the applications that will be dealt with in the next few sections are as follows:

- The law of large numbers.
- Universal source coding.
- Sanov's theorem.
- Stein's lemma and hypothesis testing.
- Conditional probability and limit theorems.

## 12.2 THE LAW OF LARGE NUMBERS

The concept of type and type classes enables us to give an alternative interpretation to the law of large numbers. In fact, it can be used as a proof of a version of the weak law in the discrete case.



The most important property of types is that there are only a polynomial number of types, and an exponential number of sequences of each type. Since the probability of each type class depends exponentially on the relative entropy distance between the type  $P$  and the distribution  $Q$ , type classes that are far from the true distribution have exponentially smaller probability.

Given an  $\epsilon > 0$ , we can define a typical set  $T_Q^\epsilon$  of sequences for the distribution  $Q^n$  as

$$T_Q^\epsilon = \{x^n : D(P_{x^n} \| Q) \leq \epsilon\}. \quad (12.63)$$

Then the probability that  $x^n$  is not typical is

$$1 - Q^n(T_Q^\epsilon) = \sum_{P : D(P \| Q) > \epsilon} Q^n(T(P)) \quad (12.64)$$

$$\leq \sum_{P : D(P \| Q) > \epsilon} 2^{-nD(P \| Q)} \quad (\text{Theorem 12.1.4}) \quad (12.65)$$

$$\leq \sum_{P : D(P \| Q) > \epsilon} 2^{-n\epsilon} \quad (12.66)$$

$$\leq (n+1)^{|X|} 2^{-n\epsilon} \quad (\text{Theorem 12.1.1}) \quad (12.67)$$

$$= 2^{-n\left(\epsilon - |X| \frac{\log(n+1)}{n}\right)}, \quad (12.68)$$

which goes to 0 as  $n \rightarrow \infty$ . Hence, the probability of the typical set goes to 1 as  $n \rightarrow \infty$ . This is similar to the AEP proved in Chapter 3, which is a form of the weak law of large numbers.

**Theorem 12.2.1:** *Let  $X_1, X_2, \dots, X_n$  be i.i.d.  $\sim P(x)$ . Then*

$$\Pr\{D(P_{x^n} \| P) > \epsilon\} \leq 2^{-n\left(\epsilon - |X| \frac{\log(n+1)}{n}\right)}, \quad (12.69)$$

*and consequently,  $D(P_{x^n} \| P) \rightarrow 0$  with probability 1.*

**Proof:** The inequality (12.69) was proved in (12.68). Summing over  $n$ , we find

$$\sum_{n=1}^{\infty} \Pr\{D(P_{x^n} \| P) > \epsilon\} < \infty. \quad (12.70)$$

Thus the expected number of occurrences of the event  $\{D(P_{x^n} \| P) > \epsilon\}$  for all  $n$  is finite, which implies that the actual number of such occurrences is also finite with probability 1 (Borel-Cantelli lemma). Hence  $D(P_{x^n} \| P) \rightarrow 0$  with probability 1.  $\square$

We now define a stronger version of typicality.

**Definition:** We will define the *strongly typical set*  $A_\epsilon^{(n)}$  to be the set of sequences in  $\mathcal{X}^n$  for which the sample frequencies are close to the true values, i.e.,

$$A_\epsilon^{(n)} = \left\{ \mathbf{x} \in \mathcal{X}^n : \left| \frac{1}{n} N(a|\mathbf{x}) - P(a) \right| < \frac{\epsilon}{|\mathcal{X}|}, \text{ for all } a \in \mathcal{X} \right\} \quad (12.71)$$

Hence the typical set consists of sequences whose type does not differ from the true probabilities by more than  $\epsilon/|\mathcal{X}|$  in any component.

By the strong law of large numbers, it follows that the probability of the strongly typical set goes to 1 as  $n \rightarrow \infty$ .

The additional power afforded by strong typicality is useful in proving stronger results, particularly in universal coding, rate distortion theory and large deviation theory.

### 12.3 UNIVERSAL SOURCE CODING

Huffman coding compresses an i.i.d. source with a known distribution  $p(x)$  to its entropy limit  $H(X)$ . However, if the code is designed for some incorrect distribution  $q(x)$ , a penalty of  $D(p||q)$  is incurred. Thus Huffman coding is sensitive to the assumed distribution.

What compression can be achieved if the true distribution  $p(x)$  is unknown? Is there a universal code of rate  $R$ , say, that suffices to describe every i.i.d. source with entropy  $H(X) < R$ ? The surprising answer is yes.

The idea is based on the method of types. There are  $2^{nH(P)}$  sequences of type  $P$ . Since there are only a polynomial number of types with denominator  $n$ , an enumeration of all sequences  $x^n$  with type  $P_{x^n}$  such that  $H(P_{x^n}) < R$  will require roughly  $nR$  bits. Thus, by describing all such sequences, we are prepared to describe any sequence that is likely to arise from any distribution  $Q$  with  $H(Q) < R$ . We begin with a definition.

**Definition:** A *fixed rate block code* of rate  $R$  for a source  $X_1, X_2, \dots, X_n$  which has an unknown distribution  $Q$  consists of two mappings, the encoder,

$$f_n: \mathcal{X}^n \rightarrow \{1, 2, \dots, 2^{nR}\}, \quad (12.72)$$

and the decoder,

$$\phi_n: \{1, 2, \dots, 2^{nR}\} \rightarrow \mathcal{X}^n. \quad (12.73)$$

Here  $R$  is called the rate of the code. The probability of error for the code with respect to the distribution  $Q$  is

$$P_e^{(n)} = Q^n(X_1, X_2, \dots, X_n : \phi_n(f_n(X_1, X_2, \dots, X_n)) \neq (X_1, X_2, \dots, X_n)). \quad (12.74)$$

**Definition:** A rate  $R$  block code for a source will be called *universal* if the functions  $f_n$  and  $\phi_n$  do not depend on the distribution  $Q$  and if  $P_e^{(n)} \rightarrow 0$  as  $n \rightarrow \infty$  if  $R > H(Q)$ .

We now describe one such universal encoding scheme, due to Csiszár and Körner [83], that is based on the fact that the number of sequences of the type  $P$  increases exponentially with the entropy and the fact that there are only a polynomial number of types.

**Theorem 12.3.1:** *There exists a sequence of  $(2^{nR}, n)$  universal source codes such that  $P_e^{(n)} \rightarrow 0$  for every source  $Q$  such that  $H(Q) < R$ .*

**Proof:** Fix the rate  $R$  for the code. Let

$$R_n = R - |\mathcal{X}| \frac{\log(n+1)}{n}. \quad (12.75)$$

Consider the set of sequences

$$A = \{\mathbf{x} \in \mathcal{X}^n : H(P_{\mathbf{x}}) \leq R_n\}. \quad (12.76)$$

Then

$$|A| = \sum_{P \in \mathcal{P}_n : H(P) \leq R_n} |T(P)| \quad (12.77)$$

$$\leq \sum_{P \in \mathcal{P}_n : H(P) \leq R_n} 2^{nH(P)} \quad (12.78)$$

$$\leq \sum_{P \in \mathcal{P}_n : H(P) \leq R_n} 2^{nR_n} \quad (12.79)$$

$$\leq (n+1)^{|\mathcal{X}|} 2^{nR_n} \quad (12.80)$$

$$= 2^{n(R_n + |\mathcal{X}| \frac{\log(n+1)}{n})} \quad (12.81)$$

$$= 2^{nR}. \quad (12.82)$$

By indexing the elements of  $A$ , we define the encoding  $f_n$  as

$$f_n(\mathbf{x}) = \begin{cases} \text{index of } \mathbf{x} \text{ in } A & \text{if } \mathbf{x} \in A, \\ 0 & \text{otherwise.} \end{cases} \quad (12.83)$$

The decoding function maps each index onto the corresponding element of  $A$ . Hence all the elements of  $A$  are encoded correctly, and all the remaining sequences result in an error. The set of sequences that are encoded correctly is illustrated in Figure 12.1.

We will now show that this encoding scheme is universal. Assume that the distribution of  $X_1, X_2, \dots, X_n$  is  $Q$  and  $H(Q) < R$ . Then the probability of decoding error is given by

$$P_e^{(n)} = 1 - Q^n(A) \quad (12.84)$$

$$= \sum_{P: H(P) > R_n} Q^n(T(P)) \quad (12.85)$$

$$\leq (n+1)^{|X|} \max_{P: H(P) > R_n} Q^n(T(P)) \quad (12.86)$$

$$\leq (n+1)^{|X|} 2^{-n \min_{P: H(P) > R_n} D(P||Q)} \quad (12.87)$$

Since  $R_n \uparrow R$  and  $H(Q) < R$ , there exists  $n_0$  such that for all  $n \geq n_0$ ,  $R_n > H(Q)$ . Then for  $n \geq n_0$ ,  $\min_{P: H(P) > R_n} D(P||Q)$  must be greater than 0, and the probability of error  $P_e^{(n)}$  converges to 0 exponentially fast as  $n \rightarrow \infty$ .

On the other hand, if the distribution  $Q$  is such that the entropy  $H(Q)$  is greater than the rate  $R$ , then with high probability, the sequence will have a type outside the set  $A$ . Hence, in such cases the probability of error is close to 1.

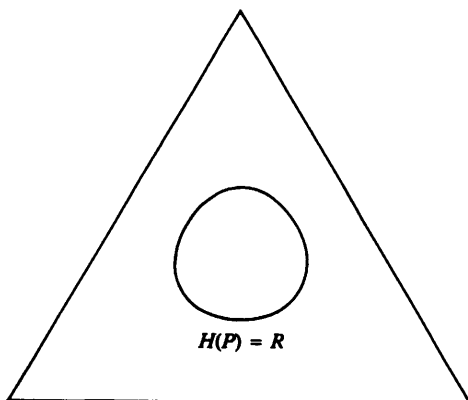


Figure 12.1. Universal code and the probability simplex.

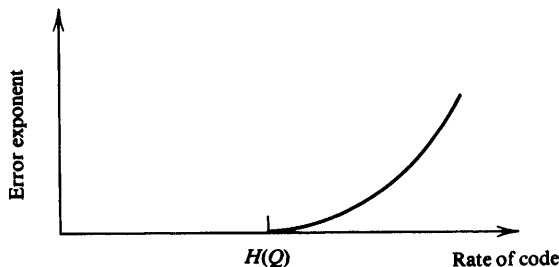


Figure 12.2. Error exponent for the universal code.

The exponent in the probability of error is

$$D_{R,Q}^* = \min_{P: H(P) > R} D(P||Q), \quad (12.88)$$

which is illustrated in Figure 12.2.  $\square$

The universal coding scheme described here is only one of many such schemes. It is universal over the set of i.i.d. distributions. There are other schemes like the Lempel-Ziv algorithm, which is a variable rate universal code for all ergodic sources. The Lempel-Ziv algorithm, discussed in Section 12.10, is often used in practice to compress data which cannot be modeled simply, such as English text or computer source code.

One may wonder why it is ever necessary to use Huffman codes, which are specific to a probability distribution. What do we lose in using a universal code?

Universal codes need a longer block length to obtain the same performance as a code designed specifically for the probability distribution. We pay the penalty for this increase in block length by the increased complexity of the encoder and decoder. Hence a distribution specific code is best if one knows the distribution of the source.

## 12.4 LARGE DEVIATION THEORY

The subject of large deviation theory can be illustrated by an example. What is the probability that  $\frac{1}{n} \sum X_i$  is near  $1/3$ , if  $X_1, X_2, \dots, X_n$  are drawn i.i.d. Bernoulli( $1/3$ )? This is a small deviation (from the expected outcome) and the probability is near 1. Now what is the probability that  $\frac{1}{n} \sum X_i$  is greater than  $3/4$  given that  $X_1, X_2, \dots, X_n$  are Bernoulli( $1/3$ )? This is a large deviation, and the probability is exponentially small. We might estimate the exponent using the central limit theorem, but this is a poor approximation for more than a few standard deviations. We note that  $\frac{1}{n} \sum X_i = 3/4$  is equivalent to  $P_{\mathbf{x}} = (1/4, 3/4)$ . Thus the probability

that  $\bar{X}_n$  is near  $3/4$  is the probability of the corresponding type. The probability of this large deviation will turn out to be  $\approx 2^{-nD((\frac{1}{2}, \frac{1}{2}) \parallel (\frac{1}{4}, \frac{3}{4}))}$ . In this section, we estimate the probability of a set of non-typical types.

Let  $E$  be a subset of the set of probability mass functions. For example,  $E$  may be the set of probability mass functions with mean  $\mu$ . With a slight abuse of notation, we write

$$Q^n(E) = Q^n(E \cap \mathcal{P}_n) = \sum_{\mathbf{x}: P_{\mathbf{x}} \in E \cap \mathcal{P}_n} Q^n(\mathbf{x}). \quad (12.89)$$

If  $E$  contains a relative entropy neighborhood of  $Q$ , then by the weak law of large numbers (Theorem 12.2.1),  $Q^n(E) \rightarrow 1$ . On the other hand, if  $E$  does not contain  $Q$  or a neighborhood of  $Q$ , then by the weak law of large numbers,  $Q^n(E) \rightarrow 0$  exponentially fast. We will use the method of types to calculate the exponent.

Let us first give some examples of the kind of sets  $E$  that we are considering. For example, assume that by observation we find that the sample average of  $g(X)$  is greater than or equal to  $\alpha$ , i.e.,  $\frac{1}{n} \sum_i g(X_i) \geq \alpha$ . This event is equivalent to the event  $P_{\mathbf{x}} \in E \cap \mathcal{P}_n$ , where

$$E = \left\{ P: \sum_{a \in \mathcal{X}} g(a)P(a) \geq \alpha \right\}, \quad (12.90)$$

because

$$\frac{1}{n} \sum_{i=1}^n g(X_i) \geq \alpha \Leftrightarrow \sum_{a \in \mathcal{X}} P_{\mathbf{x}}(a)g(a) \geq \alpha \quad (12.91)$$

$$\Leftrightarrow P_{\mathbf{x}} \in E \cap \mathcal{P}_n. \quad (12.92)$$

Thus

$$\Pr\left(\frac{1}{n} \sum_{i=1}^n g(X_i) \geq \alpha\right) = Q^n(E \cap \mathcal{P}_n) = Q^n(E). \quad (12.93)$$

Here  $E$  is a half space in the space of probability vectors, as illustrated in Figure 12.3.

**Theorem 12.4.1** (*Sanov's theorem*): Let  $X_1, X_2, \dots, X_n$  be i.i.d.  $\sim Q(x)$ . Let  $E \subseteq \mathcal{P}$  be a set of probability distributions. Then

$$Q^n(E) = Q^n(E \cap \mathcal{P}_n) \leq (n+1)^{|\mathcal{X}|} 2^{-nD(P^* \parallel Q)}, \quad (12.94)$$

where

$$P^* = \arg \min_{P \in E} D(P \parallel Q), \quad (12.95)$$

is the distribution in  $E$  that is closest to  $Q$  in relative entropy.

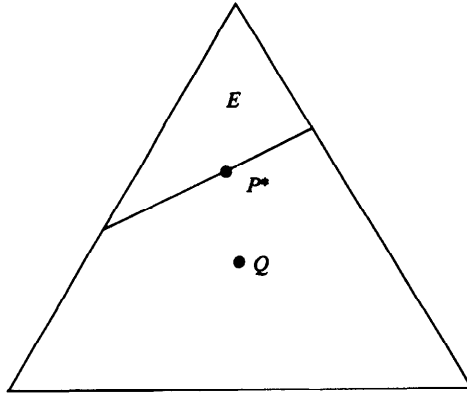


Figure 12.3. The probability simplex and Sanov's theorem.

If, in addition, the set  $E$  is the closure of its interior, then

$$\frac{1}{n} \log Q^n(E) \rightarrow -D(P^* \| Q). \quad (12.96)$$

**Proof:** We first prove the upper bound:

$$Q^n(E) = \sum_{P \in E \cap \mathcal{P}_n} Q^n(T(P)) \quad (12.97)$$

$$\leq \sum_{P \in E \cap \mathcal{P}_n} 2^{-nD(P \| Q)} \quad (12.98)$$

$$\leq \sum_{P \in E \cap \mathcal{P}_n} \max_{P \in E \cap \mathcal{P}_n} 2^{-nD(P \| Q)} \quad (12.99)$$

$$= \sum_{P \in E \cap \mathcal{P}_n} 2^{-n \min_{P \in E \cap \mathcal{P}_n} D(P \| Q)} \quad (12.100)$$

$$\leq \sum_{P \in E \cap \mathcal{P}_n} 2^{-n \min_{P \in E} D(P \| Q)} \quad (12.101)$$

$$= \sum_{P \in E \cap \mathcal{P}_n} 2^{-nD(P^* \| Q)} \quad (12.102)$$

$$\leq (n+1)^{|\mathcal{X}|} 2^{-nD(P^* \| Q)} \quad (12.103)$$

where the last inequality follows from Theorem 12.1.1.

Note that  $P^*$  need not be a member of  $\mathcal{P}_n$ . We now come to the lower bound, for which we need a “nice” set  $E$ , so that for all large  $n$ , we can find a distribution in  $E \cap \mathcal{P}_n$  which is close to  $P^*$ . If we now assume that  $E$  is the closure of its interior (thus the interior must be non-empty),

then since  $\bigcup_n \mathcal{P}_n$  is dense in the set of all distributions, it follows that  $E \cap \mathcal{P}_n$  is non-empty for all  $n \geq n_0$  for some  $n_0$ . We can then find a sequence of distributions  $P_n$  such that  $P_n \in E \cap \mathcal{P}_n$  and  $D(P_n \| Q) \rightarrow D(P^* \| Q)$ . For each  $n \geq n_0$ ,

$$Q^n(E) = \sum_{P \in E \cap \mathcal{P}_n} Q^n(T(P)) \quad (12.104)$$

$$\geq Q^n(T(P_n)) \quad (12.105)$$

$$\geq \frac{1}{(n+1)^{|\mathcal{X}|}} 2^{-nD(P_n \| Q)}. \quad (12.106)$$

Consequently,

$$\begin{aligned} \liminf \frac{1}{n} \log Q^n(E) &\geq \liminf \left( -\frac{|\mathcal{X}| \log(n+1)}{n} - D(P_n \| Q) \right) \\ &= -D(P^* \| Q). \end{aligned} \quad (12.107)$$

Combining this with the upper bound establishes the theorem.  $\square$

This argument can also be extended to continuous distributions using quantization.

## 12.5 EXAMPLES OF SANOV'S THEOREM

Suppose we wish to find  $\Pr\{\frac{1}{n} \sum_{i=1}^n g_j(X_i) \geq \alpha_j, j = 1, 2, \dots, k\}$ . Then the set  $E$  is defined as

$$E = \left\{ P : \sum_a P(a) g_j(a) \geq \alpha_j, j = 1, 2, \dots, k \right\}. \quad (12.108)$$

To find the closest distribution in  $E$  to  $Q$ , we minimize  $D(P \| Q)$  subject to the constraints in (12.108). Using Lagrange multipliers, we construct the functional

$$J(P) = \sum_x P(x) \log \frac{P(x)}{Q(x)} + \sum_i \lambda_i \sum_x P(x) g_i(x) + \nu \sum_x P(x). \quad (12.109)$$

We then differentiate and calculate the closest distribution to  $Q$  to be of the form

$$P^*(x) = \frac{Q(x) e^{\sum_i \lambda_i g_i(x)}}{\sum_{a \in \mathcal{X}} Q(a) e^{\sum_i \lambda_i g_i(a)}}, \quad (12.110)$$



where the constants  $\lambda_i$  are chosen to satisfy the constraints. Note that if  $Q$  is uniform, then  $P^*$  is the maximum entropy distribution. Verification that  $P^*$  is indeed the minimum follows from the same kind of arguments as given in Chapter 11.

Let us consider some specific examples:

**Example 12.5.1 (Dice):** Suppose we toss a fair die  $n$  times; what is the probability that the average of the throws is greater than or equal to 4? From Sanov's theorem, it follows that

$$Q^n(E) \doteq 2^{-nD(P^*\|Q)}, \quad (12.111)$$

where  $P^*$  minimizes  $D(P\|Q)$  over all distributions  $P$  that satisfy

$$\sum_{i=1}^6 iP(i) \geq 4. \quad (12.112)$$

From (12.110), it follows that  $P^*$  has the form

$$P^*(x) = \frac{2^{\lambda x}}{\sum_{i=1}^6 2^{\lambda i}}, \quad (12.113)$$

with  $\lambda$  chosen so that  $\sum iP^*(i) = 4$ . Solving numerically, we obtain  $\lambda = 0.2519$ , and  $P^* = (0.1031, 0.1227, 0.1461, 0.1740, 0.2072, 0.2468)$ , and therefore  $D(P^*\|Q) = 0.0624$  bits. Thus, the probability that the average of 10000 throws is greater than or equal to 4 is  $\approx 2^{-624}$ .

**Example 12.5.2 (Coins):** Suppose we have a fair coin, and want to estimate the probability of observing more than 700 heads in a series of 1000 tosses. The problem is like the previous example. The probability is

$$P(\bar{X}_n \geq 0.7) \doteq 2^{-nD(P^*\|Q)} \quad (12.114)$$

where  $P^*$  is the  $(0.7, 0.3)$  distribution and  $Q$  is the  $(0.5, 0.5)$  distribution. In this case,  $D(P^*\|Q) = 1 - H(P^*) = 1 - H(0.7) = 0.119$ . Thus the probability of 700 or more heads in 1000 trials is approximately  $2^{-119}$ .

**Example 12.5.3 (Mutual dependence):** Let  $Q(x, y)$  be a given joint distribution and let  $Q_0(x, y) = Q(x)Q(y)$  be the associated product distribution formed from the marginals of  $Q$ . We wish to know the likelihood that a sample drawn according to  $Q_0$  will "appear" to be jointly distributed according to  $Q$ . Accordingly, let  $(X_i, Y_i)$  be i.i.d.  $\sim Q_0(x, y) = Q(x)Q(y)$ . We define joint typicality as we did in Section 8.6,

i.e.,  $(x^n, y^n)$  is jointly typical with respect to a joint distribution  $Q(x, y)$  iff the sample entropies are close to their true values, i.e.,

$$\left| -\frac{1}{n} \log Q(x^n) - H(X) \right| \leq \epsilon, \quad (12.115)$$

$$\left| -\frac{1}{n} \log Q(y^n) - H(Y) \right| \leq \epsilon, \quad (12.116)$$

and

$$\left| -\frac{1}{n} \log Q(x^n, y^n) - H(X, Y) \right| \leq \epsilon. \quad (12.117)$$

We wish to calculate the probability (under the product distribution) of seeing a pair  $(x^n, y^n)$  that looks jointly typical of  $Q$ , i.e.,  $(x^n, y^n)$  satisfies (12.115)–(12.117). Thus  $(x^n, y^n)$  are jointly typical with respect to  $Q(x, y)$  if  $P_{x^n, y^n} \in E \cap \mathcal{P}_n(X, Y)$ , where

$$\begin{aligned} E = \left\{ P(x, y) : \left| -\sum_{x, y} P(x, y) \log Q(x) - H(X) \right| \leq \epsilon, \right. \\ \left| -\sum_{x, y} P(x, y) \log Q(y) - H(Y) \right| \leq \epsilon, \\ \left. \left| -\sum_{x, y} P(x, y) \log Q(x, y) - H(X, Y) \right| \leq \epsilon \right\}. \end{aligned} \quad (12.118)$$

Using Sanov's theorem, the probability is

$$Q_0^n(E) \doteq 2^{-nD(P^* \| Q_0)}, \quad (12.119)$$

where  $P^*$  is the distribution satisfying the constraints that is closest to  $Q_0$  in relative entropy. In this case, as  $\epsilon \rightarrow 0$ , it can be verified (Problem 10) that  $P^*$  is the joint distribution  $Q$ , and  $Q_0$  is the product distribution, so that the probability is  $2^{-nD(Q(x, y) \| Q(x)Q(y))} = 2^{-nI(X; Y)}$ , which is the same as the result derived in Chapter 8 for the joint AEP.

In the next section, we consider the empirical distribution of the sequence of outcomes given that the type is in a particular set of distributions  $E$ . We will show that not only is the probability of the set  $E$  essentially determined by  $D(P^* \| Q)$ , the distance of the closest element of  $E$  to  $Q$ , but also that the conditional type is essentially  $P^*$ , so that given that we are in set  $E$ , the type is very likely to be close to  $P^*$ .

## 12.6 THE CONDITIONAL LIMIT THEOREM

It has been shown that the probability of a set of types under a distribution  $Q$  is essentially determined by the probability of the closest element of the set to  $Q$ ; the probability is  $2^{-nD^*}$  to first order in the exponent, where

$$D^* = \min_{P \in E} D(P \| Q). \quad (12.120)$$

This follows because the probability of the set of types is the sum of the probabilities of each type, which is bounded by the largest term times the number of terms. Since the number of terms is polynomial in the length of the sequences, the sum is equal to the largest term to first order in the exponent.

We now strengthen the argument to show that not only is the probability of the set  $E$  essentially the same as the probability of the closest type  $P^*$  but also that the total probability of other types that are far away from  $P^*$  is negligible. This implies that with very high probability, the observed type is close to  $P^*$ . We call this a conditional limit theorem.

Before we prove this result, we prove a “Pythagorean” theorem, which gives some insight into the geometry of  $D(P \| Q)$ . Since  $D(P \| Q)$  is not a metric, many of the intuitive properties of distance are not valid for  $D(P \| Q)$ . The next theorem shows a sense in which  $D(P \| Q)$  behaves like the square of the Euclidean metric (Figure 12.4).

**Theorem 12.6.1:** *For a closed convex set  $E \subset \mathcal{P}$  and distribution  $Q \notin E$ , let  $P^* \in E$  be the distribution that achieves the minimum distance to  $Q$ , i.e.,*

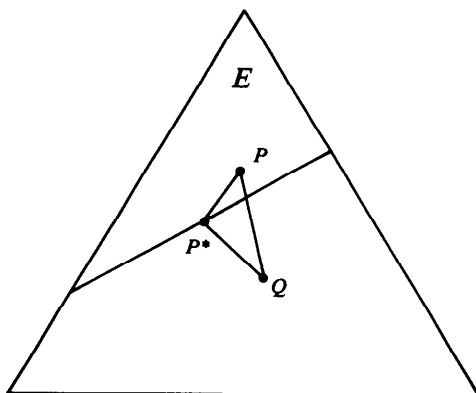


Figure 12.4. Pythagorean theorem for relative entropy.

$$D(P^* \| Q) = \min_{P \in E} D(P \| Q). \quad (12.121)$$

Then

$$D(P \| Q) \geq D(P \| P^*) + D(P^* \| Q) \quad (12.122)$$

for all  $P \in E$ .

**Note:** The main use of this theorem is as follows: suppose we have a sequence  $P_n \in E$  that yields  $D(P_n \| Q) \rightarrow D(P^* \| Q)$ . Then from the Pythagorean theorem,  $D(P_n \| P^*) \rightarrow 0$  as well.

**Proof:** Consider any  $P \in E$ . Let

$$P_\lambda = \lambda P + (1 - \lambda)P^*. \quad (12.123)$$

Then  $P_\lambda \rightarrow P^*$  as  $\lambda \rightarrow 0$ . Also since  $E$  is convex,  $P_\lambda \in E$  for  $0 \leq \lambda \leq 1$ . Since  $D(P^* \| Q)$  is the minimum of  $D(P_\lambda \| Q)$  along the path  $P^* \rightarrow P$ , the derivative of  $D(P_\lambda \| Q)$  as a function of  $\lambda$  is non-negative at  $\lambda = 0$ . Now

$$D_\lambda = D(P_\lambda \| Q) = \sum P_\lambda(x) \log \frac{P_\lambda(x)}{Q(x)}, \quad (12.124)$$

and

$$\frac{dD_\lambda}{d\lambda} = \sum \left( (P(x) - P^*(x)) \log \frac{P_\lambda(x)}{Q(x)} + (P(x) - P^*(x)) \right). \quad (12.125)$$

Setting  $\lambda = 0$ , so that  $P_\lambda = P^*$  and using the fact that  $\sum P(x) = \sum P^*(x) = 1$ , we have

$$0 \leq \left( \frac{dD_\lambda}{d\lambda} \right)_{\lambda=0} \quad (12.126)$$

$$= \sum (P(x) - P^*(x)) \log \frac{P^*(x)}{Q(x)} \quad (12.127)$$

$$= \sum P(x) \log \frac{P^*(x)}{Q(x)} - \sum P^*(x) \log \frac{P^*(x)}{Q(x)} \quad (12.128)$$

$$= \sum P(x) \log \frac{P(x)}{Q(x)} \frac{P^*(x)}{P(x)} - \sum P^*(x) \log \frac{P^*(x)}{Q(x)} \quad (12.129)$$

$$= D(P \| Q) - D(P \| P^*) - D(P^* \| Q), \quad (12.130)$$

which proves the theorem.  $\square$

Note that the relative entropy  $D(P\|Q)$  behaves like the square of the Euclidean distance. Suppose we have a convex set  $E$  in  $\mathcal{R}^n$ . Let  $A$  be a point outside the set,  $B$  the point in the set closest to  $A$ , and  $C$  any other point in the set. Then the angle between the lines  $BA$  and  $BC$  must be obtuse, which implies that  $l_{AC}^2 \geq l_{AB}^2 + l_{BC}^2$ , which is of the same form as the above theorem. This is illustrated in Figure 12.5.

We now prove a useful lemma which shows that convergence in relative entropy implies convergence in the  $\mathcal{L}_1$  norm.

**Definition:** The  $\mathcal{L}_1$  distance between any two distributions is defined as

$$\|P_1 - P_2\|_1 = \sum_{a \in \mathcal{X}} |P_1(a) - P_2(a)|. \quad (12.131)$$

Let  $A$  be the set on which  $P_1(x) > P_2(x)$ . Then

$$\|P_1 - P_2\|_1 = \sum_{x \in \mathcal{X}} |P_1(x) - P_2(x)| \quad (12.132)$$

$$= \sum_{x \in A} (P_1(x) - P_2(x)) + \sum_{x \in A^c} (P_2(x) - P_1(x)) \quad (12.133)$$

$$= P_1(A) - P_2(A) + P_2(A^c) - P_1(A^c) \quad (12.134)$$

$$= P_1(A) - P_2(A) + 1 - P_2(A) - 1 + P_1(A) \quad (12.135)$$

$$= 2(P_1(A) - P_2(A)). \quad (12.136)$$

Also note that

$$\max_{B \subseteq \mathcal{X}} (P_1(B) - P_2(B)) = P_1(A) - P_2(A) = \frac{\|P_1 - P_2\|_1}{2}. \quad (12.137)$$

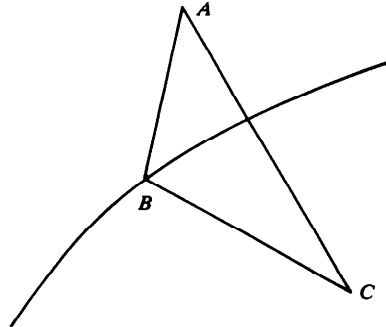


Figure 12.5. Triangle inequality for distance squared.

The left hand side of (12.137) is called the *variational distance* between  $P_1$  and  $P_2$ .

**Lemma 12.6.1:**

$$D(P_1 \| P_2) \geq \frac{1}{2 \ln 2} \|P_1 - P_2\|_1^2 \quad (12.138)$$

**Proof:** We first prove it for the binary case. Consider two binary distributions with parameters  $p$  and  $q$  with  $p \geq q$ . We will show that

$$p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q} \geq \frac{4}{2 \ln 2} (p-q)^2. \quad (12.139)$$

The difference  $g(p, q)$  between the two sides is

$$g(p, q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q} - \frac{4}{2 \ln 2} (p-q)^2. \quad (12.140)$$

Then

$$\frac{dg(p, q)}{dq} = -\frac{p}{q \ln 2} + \frac{1-p}{(1-q) \ln 2} - \frac{4}{2 \ln 2} 2(q-p) \quad (12.141)$$

$$= \frac{q-p}{q(1-q) \ln 2} - \frac{4}{\ln 2} (q-p) \quad (12.142)$$

$$\leq 0, \quad (12.143)$$

since  $q(1-q) \leq \frac{1}{4}$  and  $q \leq p$ . For  $q = p$ ,  $g(p, q) = 0$ , and hence  $g(p, q) \geq 0$  for  $q \leq p$ , which proves the lemma for the binary case.

For the general case, for any two distributions  $P_1$  and  $P_2$ , let

$$A = \{x : P_1(x) > P_2(x)\}. \quad (12.144)$$

Define a new binary random variable  $Y = \phi(X)$ , the indicator of the set  $A$ , and let  $\hat{P}_1$  and  $\hat{P}_2$  be the distributions of  $Y$ . Thus  $\hat{P}_1$  and  $\hat{P}_2$  correspond to the quantized versions of  $P_1$  and  $P_2$ . Then by the data processing inequality applied to relative entropies (which is proved in the same way as the data processing inequality for mutual information), we have

$$D(P_1 \| P_2) \geq D(\hat{P}_1 \| \hat{P}_2) \quad (12.145)$$

$$\geq \frac{4}{2 \ln 2} (P_1(A) - P_2(A))^2 \quad (12.146)$$

$$= \frac{1}{2 \ln 2} \|P_1 - P_2\|_1^2 \quad (12.147)$$

by (12.137), and the lemma is proved.  $\square$

We can now begin the proof of the conditional limit theorem. We first outline the method used. As stated at the beginning of the chapter, the essential idea is that the probability of a type under  $Q$  depends exponentially on the distance of the type from  $Q$ , and hence types that are further away are exponentially less likely to occur. We divide the set of types in  $E$  into two categories: those at about the same distance from  $Q$  as  $P^*$  and those a distance  $2\delta$  farther away. The second set has exponentially less probability than the first, and hence the first set has a conditional probability tending to 1. We then use the Pythagorean theorem to establish that all the elements in the first set are close to  $P^*$ , which will establish the theorem.

The following theorem is an important strengthening of the maximum entropy principle.

**Theorem 12.6.2** (Conditional limit theorem): *Let  $E$  be a closed convex subset of  $\mathcal{P}$  and let  $Q$  be a distribution not in  $E$ . Let  $X_1, X_2, \dots, X_n$  be discrete random variables drawn i.i.d.  $\sim Q$ . Let  $P^*$  achieve  $\min_{P \in E} D(P\|Q)$ . Then*

$$\Pr(X_1 = a | P_{X^n} \in E) \rightarrow P^*(a) \quad (12.148)$$

*in probability as  $n \rightarrow \infty$ , i.e., the conditional distribution of  $X_1$ , given that the type of the sequence is in  $E$ , is close to  $P^*$  for large  $n$ .*

**Example 12.6.1:** If  $X_i$  i.i.d.  $\sim Q$ , then

$$\Pr\left\{X_1 = a \mid \frac{1}{n} \sum X_i^2 \geq \alpha\right\} \rightarrow P^*(a), \quad (12.149)$$

where  $P^*(a)$  minimizes  $D(P\|Q)$  over  $P$  satisfying  $\sum P(a)a^2 \geq \alpha$ . This minimization results in

$$P^*(a) = Q(a) \frac{e^{\lambda a^2}}{\sum_a Q(a) e^{\lambda a^2}}, \quad (12.150)$$

where  $\lambda$  is chosen to satisfy  $\sum P^*(a)a^2 = \alpha$ . Thus the conditional distribution on  $X_1$  given a constraint on the sum of the squares is a (normalized) product of the original probability mass function and the maximum entropy probability mass function (which in this case is Gaussian).

**Proof of Theorem:** Define the sets

$$S_t = \{P \in \mathcal{P} : D(P\|Q) \leq t\}. \quad (12.151)$$

The sets  $S_t$  are convex since  $D(P\|Q)$  is a convex function of  $P$ . Let

$$D^* = D(P^*\|Q) = \min_{P \in E} D(P\|Q). \quad (12.152)$$

Then  $P^*$  is unique, since  $D(P\|Q)$  is strictly convex in  $P$ .

Now define the set

$$A = S_{D^*+2\delta} \cap E \quad (12.153)$$

and

$$B = E - S_{D^*+2\delta} \cap E. \quad (12.154)$$

Thus  $A \cup B = E$ . These sets are illustrated in Figure 12.6. Then

$$Q^n(B) = \sum_{P \in E \cap \mathcal{P}_n : D(P\|Q) > D^*+2\delta} Q^n(T(P)) \quad (12.155)$$

$$\leq \sum_{P \in E \cap \mathcal{P}_n : D(P\|Q) > D^*+2\delta} 2^{-nD(P\|Q)} \quad (12.156)$$

$$\leq \sum_{P \in E \cap \mathcal{P}_n : D(P\|Q) > D^*+2\delta} 2^{-n(D^*+2\delta)} \quad (12.157)$$

$$\leq (n+1)^{|\mathcal{X}|} 2^{-n(D^*+2\delta)}, \quad (12.158)$$

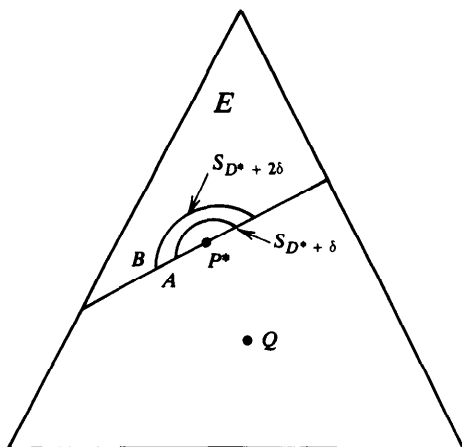


Figure 12.6. The conditional limit theorem.



since there are only a polynomial number of types. On the other hand,

$$Q^n(A) \geq Q^n(S_{D^*+\delta} \cap E) \quad (12.159)$$

$$= \sum_{P \in E \cap \mathcal{P}_n : D(P\|Q) \leq D^*+\delta} Q^n(T(P)) \quad (12.160)$$

$$\geq \sum_{P \in E \cap \mathcal{P}_n : D(P\|Q) \leq D^*+\delta} \frac{1}{(n+1)^{|\mathcal{X}|}} 2^{-nD(P\|Q)} \quad (12.161)$$

$$\geq \frac{1}{(n+1)^{|\mathcal{X}|}} 2^{-n(D^*+\delta)}, \quad \text{for } n \text{ sufficiently large,} \quad (12.162)$$

since the sum is greater than one of the terms, and for sufficiently large  $n$ , there exists at least one type in  $S_{D^*+\delta} \cap E \cap \mathcal{P}_n$ . Then for  $n$  sufficiently large

$$\Pr(P_{X^n} \in B | P_{X^n} \in E) = \frac{Q^n(B \cap E)}{Q^n(E)} \quad (12.163)$$

$$\leq \frac{Q^n(B)}{Q^n(A)} \quad (12.164)$$

$$\leq \frac{(n+1)^{|\mathcal{X}|} 2^{-n(D^*+2\delta)}}{\frac{1}{(n+1)^{|\mathcal{X}|}} 2^{-n(D^*+\delta)}} \quad (12.165)$$

$$= (n+1)^{2|\mathcal{X}|} 2^{-n\delta}, \quad (12.166)$$

which goes to 0 as  $n \rightarrow \infty$ . Hence the conditional probability of  $B$  goes to 0 as  $n \rightarrow \infty$ , which implies that the conditional probability of  $A$  goes to 1.

We now show that all the members of  $A$  are close to  $P^*$  in relative entropy. For all members of  $A$ ,

$$D(P\|Q) \leq D^* + 2\delta. \quad (12.167)$$

Hence by the “Pythagorean” theorem (Theorem 12.6.1),

$$D(P\|P^*) + D(P^*\|Q) \leq D(P\|Q) \leq D^* + 2\delta, \quad (12.168)$$

which in turn implies that

$$D(P\|P^*) \leq 2\delta, \quad (12.169)$$

since  $D(P^*\|Q) = D^*$ . Thus  $P_{\mathbf{x}} \in A$  implies that  $D(P_{\mathbf{x}}\|Q) \leq D^* + 2\delta$ , and therefore that  $D(P_{\mathbf{x}}\|P^*) \leq 2\delta$ . Consequently, since  $\Pr\{P_{X^n} \in A | P_{X^n} \in E\} \rightarrow 1$ , it follows that

$$\Pr(D(P_{X^n}||P^*) \leq 2\delta | P_{X^n} \in E) \rightarrow 1 \quad (12.170)$$

as  $n \rightarrow \infty$ .

By Lemma 12.6.1, the fact that the relative entropy is small implies that the  $\mathcal{L}_1$  distance is small, which in turn implies that  $\max_{a \in \mathcal{X}} |P_{X^n}(a) - P^*(a)|$  is small. Thus  $\Pr(|P_{X^n}(a) - P^*(a)| \geq \epsilon | P_{X^n} \in E) \rightarrow 0$  as  $n \rightarrow \infty$ . Alternatively, this can be written as

$$\Pr(X_1 = a | P_{X^n} \in E) \rightarrow P^*(a) \text{ in probability.} \quad (12.171)$$

In this theorem, we have only proved that the marginal distribution goes to  $P^*$  as  $n \rightarrow \infty$ . Using a similar argument, we can prove a stronger version of this theorem, i.e.,

$$\Pr(X_1 = a_1, X_2 = a_2, \dots, X_m = a_m | P_{X^n} \in E) \rightarrow \prod_{i=1}^m P^*(a_i) \text{ in probability.} \quad (12.172)$$

This holds for fixed  $m$  as  $n \rightarrow \infty$ . The result is not true for  $m = n$ , since there are end effects; given that the type of the sequence is in  $E$ , the last elements of the sequence can be determined from the remaining elements, and the elements are no longer independent. The conditional limit theorem states that the first few elements are asymptotically independent with common distribution  $P^*$ .

**Example 12.6.2:** As an example of the conditional limit theorem, let us consider the case when  $n$  fair dice are rolled. Suppose that the sum of the outcomes exceeds  $4n$ . Then by the conditional limit theorem, the probability that the first die shows a number  $a \in \{1, 2, \dots, 6\}$  is approximately  $P^*(a)$ , where  $P^*(a)$  is the distribution in  $E$  that is closest to the uniform distribution, where  $E = \{P: \sum P(a)a \geq 4\}$ . This is the maximum entropy distribution given by

$$P^*(x) = \frac{2^{\lambda x}}{\sum_{i=1}^6 2^{\lambda i}}, \quad (12.173)$$

with  $\lambda$  chosen so that  $\sum iP^*(i) = 4$  (see Chapter 11). Here  $P^*$  is the conditional distribution on the first (or any other) die. Apparently the first few dice inspected will behave as if they are independently drawn according to an exponential distribution.

## 12.7 HYPOTHESIS TESTING

One of the standard problems in statistics is to decide between two alternative explanations for the observed data. For example, in medical

testing, one may wish to test whether a new drug is effective or not. Similarly, a sequence of coin tosses may reveal whether the coin is biased or not.

These problems are examples of the general hypothesis testing problem. In the simplest case, we have to decide between two i.i.d. distributions. The general problem can be stated as follows:

**Problem:** Let  $X_1, X_2, \dots, X_n$  be i.i.d.  $\sim Q(x)$ . We consider two hypotheses:

- $H_1: Q = P_1$ .
- $H_2: Q = P_2$ .

Consider the general decision function  $g(x_1, x_2, \dots, x_n)$ , where  $g(x_1, x_2, \dots, x_n) = 1$  implies that  $H_1$  is accepted and  $g(x_1, x_2, \dots, x_n) = 2$  implies that  $H_2$  is accepted. Since the function takes on only two values, the test can also be specified by specifying the set  $A$  over which  $g(x_1, x_2, \dots, x_n)$  is 1; the complement of this set is the set where  $g(x_1, x_2, \dots, x_n)$  has the value 2. We define the two probabilities of error:

$$\alpha = \Pr(g(X_1, X_2, \dots, X_n) = 2 | H_1 \text{ true}) = P_1^n(A^c) \quad (12.174)$$

and

$$\beta = \Pr(g(X_1, X_2, \dots, X_n) = 1 | H_2 \text{ true}) = P_2^n(A). \quad (12.175)$$

In general, we wish to minimize both probabilities, but there is a trade-off. Thus we minimize one of the probabilities of error subject to a constraint on the other probability of error. The best achievable error exponent in the probability of error for this problem is given by Stein's lemma.

We first prove the Neyman-Pearson lemma, which derives the form of the optimum test between two hypotheses. We derive the result for discrete distributions; the same results can be derived for continuous distributions as well.

**Theorem 12.7.1 (Neyman-Pearson lemma):** Let  $X_1, X_2, \dots, X_n$  be drawn i.i.d. according to probability mass function  $P$ . Consider the decision problem corresponding to hypotheses  $Q = P_1$  vs.  $Q = P_2$ . For  $T \geq 0$ , define a region

$$A_n(T) = \left\{ \frac{P_1(x_1, x_2, \dots, x_n)}{P_2(x_1, x_2, \dots, x_n)} > T \right\}. \quad (12.176)$$

Let

$$\alpha^* = P_1^n(A_n^c(T)), \quad \beta^* = P_2^n(A_n(T)), \quad (12.177)$$

be the corresponding probabilities of error corresponding to decision region  $A_n$ . Let  $B_n$  be any other decision region with associated probabilities of error  $\alpha$  and  $\beta$ . If  $\alpha \leq \alpha^*$ , then  $\beta \geq \beta^*$ .

**Proof:** Let  $A = A_n(T)$  be the region defined in (12.176) and let  $B \in \mathcal{X}^n$  be any other acceptance region. Let  $\phi_A$  and  $\phi_B$  be the indicator functions of the decision regions  $A$  and  $B$  respectively. Then for all  $\mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathcal{X}^n$ ,

$$[\phi_A(\mathbf{x}) - \phi_B(\mathbf{x})][P_1(\mathbf{x}) - TP_2(\mathbf{x})] \geq 0. \quad (12.178)$$

This can be seen by considering separately the cases  $\mathbf{x} \in A$  and  $\mathbf{x} \notin A$ .

Multiplying out and integrating this over the entire space, we obtain

$$0 \leq \sum (\phi_A P_1 - T \phi_A P_2 - P_1 \phi_B + T P_2 \phi_B) \quad (12.179)$$

$$= \sum_A (P_1 - T P_2) - \sum_B (P_1 - T P_2) \quad (12.180)$$

$$= (1 - \alpha^*) - T \beta^* - (1 - \alpha) + T \beta \quad (12.181)$$

$$= T(\beta - \beta^*) - (\alpha^* - \alpha). \quad (12.182)$$

Since  $T \geq 0$ , we have proved the theorem.  $\square$

The Neyman-Pearson lemma indicates that the optimum test for two hypotheses is of the form

$$\frac{P_1(X_1, X_2, \dots, X_n)}{P_2(X_1, X_2, \dots, X_n)} > T. \quad (12.183)$$

This is the likelihood ratio test and the quantity  $\frac{P_1(X_1, X_2, \dots, X_n)}{P_2(X_1, X_2, \dots, X_n)}$  is called the *likelihood ratio*.

For example, in a test between two Gaussian distributions, i.e., between  $f_1 = \mathcal{N}(1, \sigma^2)$  and  $f_2 = \mathcal{N}(-1, \sigma^2)$ , the likelihood ratio becomes

$$\frac{f_1(X_1, X_2, \dots, X_n)}{f_2(X_1, X_2, \dots, X_n)} = \frac{\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(X_i-1)^2}{2\sigma^2}}}{\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(X_i+1)^2}{2\sigma^2}}} \quad (12.184)$$

$$= e^{+\frac{2 \sum_{i=1}^n X_i}{\sigma^2}} \quad (12.185)$$

$$= e^{+\frac{2n\bar{X}_n}{\sigma^2}}. \quad (12.186)$$

Hence the likelihood ratio test consists of comparing the sample mean

$\bar{X}_n$  with a threshold. If we want the two probabilities of error to be equal, we should set  $T = 1$ . This is illustrated in Figure 12.7.

In the above theorem, we have shown that the optimum test is a likelihood ratio test. We can rewrite the log-likelihood ratio as

$$L(X_1, X_2, \dots, X_n) = \log \frac{P_1(X_1, X_2, \dots, X_n)}{P_2(X_1, X_2, \dots, X_n)} \quad (12.187)$$

$$= \sum_{i=1}^n \log \frac{P_1(X_i)}{P_2(X_i)} \quad (12.188)$$

$$= \sum_{a \in \mathcal{X}} n P_{X^n}(a) \log \frac{P_1(a)}{P_2(a)} \quad (12.189)$$

$$= \sum_{a \in \mathcal{X}} n P_{X^n}(a) \log \frac{P_1(a)}{P_2(a)} \frac{P_{X^n}(a)}{P_{X^n}(a)} \quad (12.190)$$

$$= \sum_{a \in \mathcal{X}} n P_{X^n}(a) \log \frac{P_{X^n}(a)}{P_2(a)} - \sum_{a \in \mathcal{X}} n P_{X^n}(a) \log \frac{P_{X^n}(a)}{P_1(a)} \quad (12.191)$$

$$= nD(P_{X^n} \| P_2) - nD(P_{X^n} \| P_1), \quad (12.192)$$

the difference between the relative entropy distances of the sample type to each of the two distributions. Hence the likelihood ratio test

$$\frac{P_1(X_1, X_2, \dots, X_n)}{P_2(X_1, X_2, \dots, X_n)} > T \quad (12.193)$$

is equivalent to

$$D(P_{X^n} \| P_2) - D(P_{X^n} \| P_1) > \frac{1}{n} \log T. \quad (12.194)$$

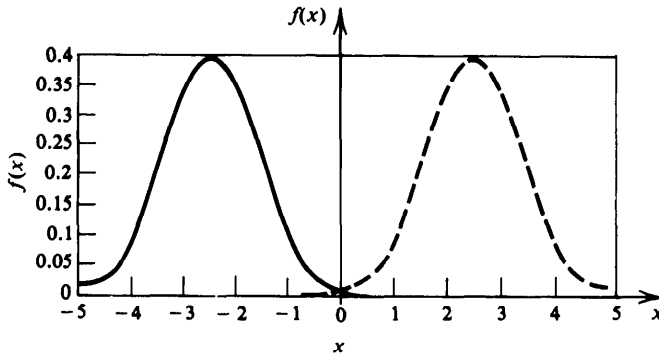


Figure 12.7. Testing between two Gaussian distributions.

We can consider the test to be equivalent to specifying a region of the simplex of types which corresponds to choosing hypothesis 1. The optimum region is of the form (12.194), for which the boundary of the region is the set of types for which the difference between the distances is a constant. This boundary is the analog of the perpendicular bisector in Euclidean geometry. The test is illustrated in Figure 12.8.

We now offer some informal arguments based on Sanov's theorem to show how to choose the threshold to obtain different probabilities of error. Let  $B$  denote the set on which hypothesis 1 is accepted. The probability of error of the first kind is

$$\alpha_n = P_1^n(P_{X^n} \in B^c). \quad (12.195)$$

Since the set  $B^c$  is convex, we can use Sanov's theorem to show that the probability of error is essentially determined by the relative entropy of the closest member of  $B^c$  to  $P_1$ . Therefore,

$$\alpha_n \doteq 2^{-nD(P_1^* \| P_1)}, \quad (12.196)$$

where  $P_1^*$  is the closest element of  $B^c$  to distribution  $P_1$ . Similarly,

$$\beta_n \doteq 2^{-nD(P_2^* \| P_2)}, \quad (12.197)$$

where  $P_2^*$  is the closest element in  $B$  to the distribution  $P_2$ .

Now minimizing  $D(P \| P_2)$  subject to the constraint  $D(P \| P_2) - D(P \| P_1) \geq \frac{1}{n} \log T$  will yield the type in  $B$  that is closest to  $P_2$ . Setting up the minimization of  $D(P \| P_2)$  subject to  $D(P \| P_2) - D(P \| P_1) = \frac{1}{n} \log T$  using Lagrange multipliers, we have

$$J(P) = \sum P(x) \log \frac{P(x)}{P_2(x)} + \lambda \sum P(x) \log \frac{P_1(x)}{P_2(x)} + \nu \sum P(x). \quad (12.198)$$

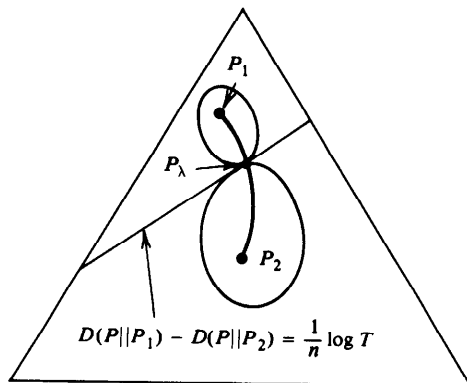


Figure 12.8. The likelihood ratio test on the probability simplex.

Differentiating with respect to  $P(x)$  and setting to 0, we have

$$\log \frac{P(x)}{P_2(x)} + 1 + \lambda \log \frac{P_1(x)}{P_2(x)} + \nu = 0. \quad (12.199)$$

Solving this set of equations, we obtain the minimizing  $P$  of the form

$$P_2^* = P_{\lambda^*} = \frac{P_1^\lambda(x) P_2^{1-\lambda}(x)}{\sum_{a \in \mathcal{X}} P_1^\lambda(a) P_2^{1-\lambda}(a)}, \quad (12.200)$$

where  $\lambda$  is chosen so that  $D(P_{\lambda^*} \| P_1) - D(P_{\lambda^*} \| P_2) = \frac{\log T}{n}$ .

From the symmetry of expression (12.200), it is clear that  $P_1^* = P_2^*$  and that the probabilities of error behave exponentially with exponents given by the relative entropies  $D(P^* \| P_1)$  and  $D(P^* \| P_2)$ . Also note from the equation that as  $\lambda \rightarrow 1$ ,  $P_\lambda \rightarrow P_1$  and as  $\lambda \rightarrow 0$ ,  $P_\lambda \rightarrow P_2$ . The line that  $P_\lambda$  traces out as  $\lambda$  varies is a geodesic in the simplex. Here  $P_\lambda$  is a normalized convex combination, where the combination is in the exponent (Figure 12.8).

In the next section, we calculate the best error exponent when one of the two types of error goes to zero arbitrarily slowly (Stein's lemma). We will also minimize the weighted sum of the two probabilities of error and obtain the Chernoff bound.

## 12.8 STEIN'S LEMMA

We now consider the case when one of the probabilities of error is fixed and we wish to minimize the other probability of error. The best error exponent in this case is given by Stein's lemma.

**Theorem 12.8.1 (Stein's lemma):** Let  $X_1, X_2, \dots, X_n$  be i.i.d.  $\sim Q$ . Consider the hypothesis test between two alternatives,  $Q = P_1$  and  $Q = P_2$ , where  $D(P_1 \| P_2) < \infty$ . Let  $A_n \subseteq \mathcal{X}^n$  be an acceptance region for hypothesis 1. Let the probabilities of error be

$$\alpha_n = P_1^n(A_n^c), \quad \beta_n = P_2^n(A_n). \quad (12.201)$$

and for  $0 < \epsilon < \frac{1}{2}$ , define

$$\beta_n^\epsilon = \min_{\substack{A_n \subseteq \mathcal{X}^n \\ \alpha_n < \epsilon}} \beta_n. \quad (12.202)$$

Then

$$\lim_{\epsilon \rightarrow 0} \lim_{n \rightarrow \infty} \frac{1}{n} \log \beta_n^\epsilon = -D(P_1 \| P_2). \quad (12.203)$$

**Proof:** To prove the theorem, we construct a sequence of acceptance regions  $A_n \subseteq \mathcal{X}^n$  such that  $\alpha_n < \epsilon$  and  $\beta_n \doteq 2^{-nD(P_1\|P_2)}$ . We then show that no other sequence of tests has an asymptotically better exponent.

First, we define

$$A_n = \left\{ \mathbf{x} \in \mathcal{X}^n : 2^{+n(D(P_1\|P_2)-\delta)} \leq \frac{P_1(\mathbf{x})}{P_2(\mathbf{x})} \leq 2^{+n(D(P_1\|P_2)+\delta)} \right\}. \quad (12.204)$$

Then we have the following properties:

1.  $P_1^n(A_n) \rightarrow 1$ . This follows from

$$P_1^n(A_n) = P_1^n \left( \frac{1}{n} \sum_{i=1}^n \log \frac{P_1(X_i)}{P_2(X_i)} \in (D(P_1\|P_2) - \delta, D(P_1\|P_2) + \delta) \right) \quad (12.205)$$

$$\rightarrow 1 \quad (12.206)$$

by the strong law of large numbers, since  $D(P_1\|P_2) = E_{P_1} \left( \log \frac{P_1(X)}{P_2(X)} \right)$ . Hence for sufficiently large  $n$ ,  $\alpha_n < \epsilon$ .

2.  $P_2^n(A_n) \leq 2^{-n(D(P_1\|P_2)-\delta)}$ . Using the definition of  $A_n$ , we have

$$P_2^n(A_n) = \sum_{A_n} P_2(\mathbf{x}) \quad (12.207)$$

$$\leq \sum_{A_n} P_1(\mathbf{x}) 2^{-n(D(P_1\|P_2)-\delta)} \quad (12.208)$$

$$= 2^{-n(D(P_1\|P_2)-\delta)} \sum_{A_n} P_1(\mathbf{x}) \quad (12.209)$$

$$= 2^{-n(D(P_1\|P_2)-\delta)} (1 - \alpha_n). \quad (12.210)$$

Similarly, we can show that

$$P_2^n(A_n) \geq 2^{-n(D(P_1\|P_2)+\delta)} (1 - \alpha_n). \quad (12.211)$$

Hence

$$\frac{1}{n} \log \beta_n \leq -D(P_1\|P_2) + \delta + \frac{\log(1 - \alpha_n)}{n}, \quad (12.212)$$

and

$$\frac{1}{n} \log \beta_n \geq -D(P_1\|P_2) - \delta + \frac{\log(1 - \alpha_n)}{n}, \quad (12.213)$$



Hence

$$\lim_{n \rightarrow \infty} \lim \frac{1}{n} \log \beta_n = -D(P_1 \| P_2). \quad (12.214)$$

3. We now prove that no other sequence of acceptance regions does better. Let  $B_n \subseteq \mathcal{X}^n$  be any other sequence of acceptance regions with  $\alpha_{n, B_n} = P_1^n(B_n^c) < \epsilon$ . Let  $\beta_{n, B_n} = P_2^n(B_n)$ . We will show that  $\beta_{n, B_n} \geq 2^{-n(D(P_1 \| P_2) - \epsilon)}$ .

Here

$$\beta_{n, B_n} = P_2^n(B_n) \geq P_2^n(A_n \cap B_n) \quad (12.215)$$

$$= \sum_{A_n \cap B_n} P_2(\mathbf{x}) \quad (12.216)$$

$$\geq \sum_{A_n \cap B_n} P_1(\mathbf{x}) 2^{-n(D(P_1 \| P_2) + \delta)} \quad (12.217)$$

$$= 2^{-n(D(P_1 \| P_2) + \delta)} \sum_{A_n \cap B_n} P_1(\mathbf{x}) \quad (12.218)$$

$$\geq (1 - \alpha_n - \alpha_{n, B_n}) 2^{-n(D(P_1 \| P_2) + \delta)}, \quad (12.219)$$

where the last inequality follows from the union of events bound as follows:

$$\sum_{A_n \cap B_n} P_1(\mathbf{x}) = P_1(A_n \cap B_n) \quad (12.220)$$

$$= 1 - P_1(A_n^c \cup B_n^c) \quad (12.221)$$

$$\geq 1 - P_1(A_n^c) - P_1(B_n^c) \quad (12.222)$$

$$= 1 - \alpha_n - \alpha_{n, B_n}. \quad (12.223)$$

Hence

$$\frac{1}{n} \log \beta_{n, B_n} \geq -D(P_1 \| P_2) - \delta - \frac{\log(1 - \alpha_n - \alpha_{n, B_n})}{n}, \quad (12.224)$$

and since  $\delta > 0$  is arbitrary,

$$\lim_{n \rightarrow \infty} \lim \frac{1}{n} \log \beta_{n, B_n} \geq -D(P_1 \| P_2). \quad (12.225)$$

Thus no sequence of sets  $B_n$  has an exponent better than  $D(P_1 \| P_2)$ . But the sequence  $A_n$  achieves the exponent  $D(P_1 \| P_2)$ . Thus  $A_n$  is asymptotically optimal, and the best error exponent is  $D(P_1 \| P_2)$ .  $\square$

## 12.9 CHERNOFF BOUND

We have considered the problem of hypothesis testing in the classical setting, in which we treat the two probabilities of error separately. In the derivation of Stein's lemma, we set  $\alpha_n \leq \epsilon$  and achieved  $\beta_n \doteq 2^{-nD}$ . But this approach lacks symmetry. Instead, we can follow a Bayesian approach, in which we assign prior probabilities to both the hypotheses. In this case, we wish to minimize the overall probability of error given by the weighted sum of the individual probabilities of error. The resulting error exponent is the Chernoff information.

The setup is as follows:  $X_1, X_2, \dots, X_n$  i.i.d.  $\sim Q$ . We have two hypotheses:  $Q = P_1$  with prior probability  $\pi_1$  and  $Q = P_2$  with prior probability  $\pi_2$ . The overall probability of error is

$$P_e^{(n)} = \pi_1 \alpha_n + \pi_2 \beta_n. \quad (12.226)$$

Let

$$D^* = \lim_{n \rightarrow \infty} \min_{A_n \subseteq \mathcal{X}^n} -\frac{1}{n} \log P_e^{(n)}. \quad (12.227)$$

**Theorem 12.9.1** (*Chernoff*): *The best achievable exponent in the Bayesian probability of error is  $D^*$ , where*

$$D^* = D(P_{\lambda^*} \| P_1) = D(P_{\lambda^*} \| P_2), \quad (12.228)$$

with

$$P_{\lambda} = \frac{P_1^{\lambda}(x) P_2^{1-\lambda}(x)}{\sum_{a \in \mathcal{X}} P_1^{\lambda}(a) P_2^{1-\lambda}(a)} \quad (12.229)$$

and  $\lambda^*$  the value of  $\lambda$  such that

$$D(P_{\lambda^*} \| P_1) = D(P_{\lambda^*} \| P_2). \quad (12.230)$$

**Proof:** The basic details of the proof were given in the previous section. We have shown that the optimum test is a likelihood ratio test, which can be considered to be of the form

$$D(P_{X^n} \| P_2) - D(P_{X^n} \| P_1) > T. \quad (12.231)$$

The test divides the probability simplex into regions corresponding to hypothesis 1 and hypothesis 2, respectively. This is illustrated in Figure 12.9.

Let  $A$  be the set of types associated with hypothesis 1. From the

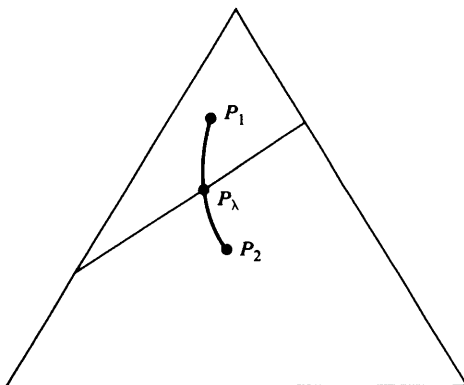


Figure 12.9. The probability simplex and Chernoff's bound.

discussion preceding (12.200), it follows that the closest point in the set  $A^c$  to  $P_1$  is on the boundary of  $A$ , and is of the form given by (12.229). Then from the discussion in the last section, it is clear that  $P_\lambda$  is the distribution in  $A$  that is closest to  $P_2$ ; it is also the distribution in  $A^c$  that is closest to  $P_1$ . By Sanov's theorem, we can calculate the associated probabilities of error

$$\alpha_n = P_1^n(A^c) \doteq 2^{-nD(P_\lambda \| P_1)} \quad (12.232)$$

and

$$\beta_n = P_2^n(A) \doteq 2^{-nD(P_\lambda \| P_2)}. \quad (12.233)$$

In the Bayesian case, the overall probability of error is the weighted sum of the two probabilities of error,

$$P_e \doteq \pi_1 2^{-nD(P_\lambda \| P_1)} + \pi_2 2^{-nD(P_\lambda \| P_2)} \doteq 2^{-n \min\{D(P_\lambda \| P_1), D(P_\lambda \| P_2)\}}, \quad (12.234)$$

since the exponential rate is determined by the worst exponent. Since  $D(P_\lambda \| P_1)$  increases with  $\lambda$  and  $D(P_\lambda \| P_2)$  decreases with  $\lambda$ , the maximum value of the minimum of  $\{D(P_\lambda \| P_1), D(P_\lambda \| P_2)\}$  is attained when they are equal. This is illustrated in Figure 12.10.

Hence, we choose  $\lambda$  so that

$$D(P_\lambda \| P_1) = D(P_\lambda \| P_2) \triangleq C(P_1, P_2). \quad (12.235)$$

Thus  $C(P_1, P_2)$  is the highest achievable exponent for the probability of error and is called the Chernoff information.  $\square$

The above definition is equivalent to the standard definition of *Chernoff information*,

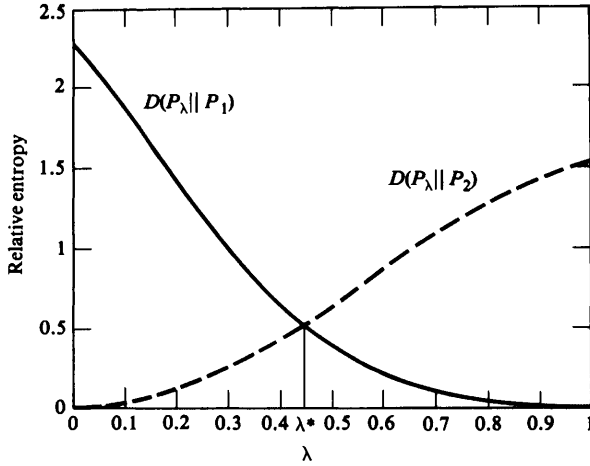


Figure 12.10. Relative entropy  $D(P_\lambda \| P_1)$  and  $D(P_\lambda \| P_2)$  as a function of  $\lambda$ .

$$C(P_1, P_2) = - \min_{0 \leq \lambda \leq 1} \log \left( \sum_x P_1^\lambda(x) P_2^{1-\lambda}(x) \right). \quad (12.236)$$

It is left as an exercise to the reader to show (algebraically) the equivalence of (12.235) and (12.236). We will briefly outline the usual derivation of the Chernoff bound. The maximum *a posteriori* probability decision rule minimizes the Bayesian probability of error. The decision region  $A$  for hypothesis 1 for the maximum *a posteriori* rule is

$$A = \left\{ \mathbf{x} : \frac{\pi_1 P_1(\mathbf{x})}{\pi_2 P_2(\mathbf{x})} > 1 \right\}, \quad (12.237)$$

the set of outcomes where the *a posteriori* probability of hypothesis 1 is greater than the *a posteriori* probability of hypothesis 2. The probability of error for this rule is

$$P_e = \pi_1 \alpha_n + \pi_2 \beta_n \quad (12.238)$$

$$= \sum_{A^c} \pi_1 P_1 + \sum_A \pi_2 P_2 \quad (12.239)$$

$$= \sum \min\{\pi_1 P_1, \pi_2 P_2\}. \quad (12.240)$$

Now for any two positive numbers  $a$  and  $b$ , we have

$$\min\{a, b\} \leq a^\lambda b^{1-\lambda}, \quad \text{for all } 0 \leq \lambda \leq 1. \quad (12.241)$$

Using this to continue the chain, we have

$$P_e = \sum \min\{\pi_1 P_1, \pi_2 P_2\} \quad (12.242)$$

$$\leq \sum (\pi_1 P_1)^\lambda (\pi_2 P_2)^{1-\lambda} \quad (12.243)$$

$$\leq \sum P_1^\lambda P_2^{1-\lambda}. \quad (12.244)$$

For a sequence of i.i.d. observations,  $P_k(\mathbf{x}) = \prod_{i=1}^n P_k(x_i)$ , and

$$P_e^{(n)} \leq \sum \pi_1^\lambda \pi_2^{1-\lambda} \prod_i P_1^\lambda(x_i) P_2^{1-\lambda}(x_i) \quad (12.245)$$

$$= \pi_1^\lambda \pi_2^{1-\lambda} \prod_i \sum P_1^\lambda(x_i) P_2^{1-\lambda}(x_i) \quad (12.246)$$

$$\stackrel{(a)}{\leq} \prod_i \sum P_1^\lambda P_2^{1-\lambda} \quad (12.247)$$

$$= \left( \sum P_1^\lambda P_2^{1-\lambda} \right)^n, \quad (12.248)$$

where (a) follows since  $\pi_1 \leq 1$ ,  $\pi_2 \leq 1$ . Hence, we have

$$\frac{1}{n} \log P_e^{(n)} \leq \log \sum P_1^\lambda(x) P_2^{1-\lambda}(x) \quad (12.249)$$

Since this is true for all  $\lambda$ , we can take the minimum over  $0 \leq \lambda \leq 1$ , resulting in the Chernoff bound. This proves that the exponent is no better than  $C(P_1, P_2)$ . Achievability follows from Theorem 12.9.1.

Note that the Bayesian error exponent does not depend on the actual value of  $\pi_1$  and  $\pi_2$ , as long as they are non-zero. Essentially, the effect of the prior is washed out for large sample sizes. The optimum decision rule is to choose the hypothesis with the maximum *a posteriori* probability, which corresponds to the test

$$\frac{\pi_1 P_1(X_1, X_2, \dots, X_n)}{\pi_2 P_2(X_1, X_2, \dots, X_n)} \geq 1. \quad (12.250)$$

Taking the log and dividing by  $n$ , this test can be rewritten as

$$\frac{1}{n} \log \frac{\pi_1}{\pi_2} + \frac{1}{n} \sum_i \log \frac{P_1(X_i)}{P_2(X_i)} \leq 0, \quad (12.251)$$

where the second term tends to  $D(P_1 \| P_2)$  or  $-D(P_2 \| P_1)$  accordingly as  $P_1$  or  $P_2$  is the true distribution. The first term tends to 0, and the effect of the prior distribution washes out.

Finally, to round off our discussion of large deviation theory and hypothesis testing, we consider an example of the conditional limit theorem.

**Example 12.9.1:** Suppose major league baseball players have a batting average of 260 with a standard deviation of 15 and suppose that minor league ballplayers have a batting average of 240 with a standard deviation of 15. A group of 100 ballplayers from one of the leagues (the league is chosen at random) is found to have a group batting average greater than 250, and is therefore judged to be major leaguers. We are now told that we are mistaken; these players are minor leaguers. What can we say about the distribution of batting averages among these 100 players? It will turn out that the distribution of batting averages among these players will have a mean of 250 and a standard deviation of 15. This follows from the conditional limit theorem. To see this, we abstract the problem as follows.

Let us consider an example of testing between two Gaussian distributions,  $f_1 = \mathcal{N}(1, \sigma^2)$  and  $f_2 = \mathcal{N}(-1, \sigma^2)$ , with different means and the same variance. As discussed in the last section, the likelihood ratio test in this case is equivalent to comparing the sample mean with a threshold. The Bayes test is "Accept the hypothesis  $f = f_1$  if  $\frac{1}{n} \sum_{i=1}^n X_i > 0$ ."

Now assume that we make an error of the first kind (we say  $f = f_1$  when indeed  $f = f_2$ ) in this test. What is the conditional distribution of the samples given that we have made an error?

We might guess at various possibilities:

- The sample will look like a  $(\frac{1}{2}, \frac{1}{2})$  mix of the two normal distributions. Plausible as this is, it is incorrect.
- $X_i \approx 0$  for all  $i$ . This is quite clearly very very unlikely, although it is conditionally likely that  $\bar{X}_n$  is close to 0.
- The correct answer is given by the conditional limit theorem. If the true distribution is  $f_2$  and the sample type is in the set  $A$ , the conditional distribution is close to  $f^*$ , the distribution in  $A$  that is closest to  $f_2$ . By symmetry, this corresponds to  $\lambda = \frac{1}{2}$  in (12.229). Calculating the distribution, we get

$$f^*(x) = \frac{\left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-1)^2}{2\sigma^2}}\right)^{1/2} \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x+1)^2}{2\sigma^2}}\right)^{1/2}}{\int \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-1)^2}{2\sigma^2}}\right)^{1/2} \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x+1)^2}{2\sigma^2}}\right)^{1/2} dx} \quad (12.252)$$

$$= \frac{\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x^2+1)}{2\sigma^2}}}{\int \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x^2+1)}{2\sigma^2}} dx} \quad (12.253)$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x^2}{2\sigma^2}} \quad (12.254)$$

$$= \mathcal{N}(0, \sigma^2). \quad (12.255)$$

It is interesting to note that the conditional distribution is normal with mean 0 and with the same variance as the original distributions. This is strange but true; if we mistake a normal population for another, the “shape” of this population still looks normal with the same variance and a different mean. Apparently, this rare event does not result from bizarre looking data.

**Example 12.9.2** (*Large deviation theory and football*): Consider a very simple version of football in which the score is directly related to the number of yards gained. Assume that the coach has a choice between two strategies: running or passing. Associated with each strategy is a distribution on the number of yards gained. For example, in general, running results in a gain of a few yards with very high probability, while passing results in huge gains with low probability. Examples of the distributions are illustrated in Figure 12.11.

At the beginning of the game, the coach uses the strategy that promises the greatest expected gain. Now assume that we are in the closing minutes of the game and one of the teams is leading by a large margin. (Let us ignore first downs and adaptable defenses.) So the trailing team will win only if it is very lucky. If luck is required to win, then we might as well assume that we will be lucky and play accordingly. What is the appropriate strategy?

Assume that the team has only  $n$  plays left and it must gain  $l$  yards, where  $l$  is much larger than  $n$  times the expected gain under each play. The probability that the team succeeds in achieving  $l$  yards is exponentially small; hence, we can use the large deviation results and Sanov’s theorem to calculate the probability of this event.

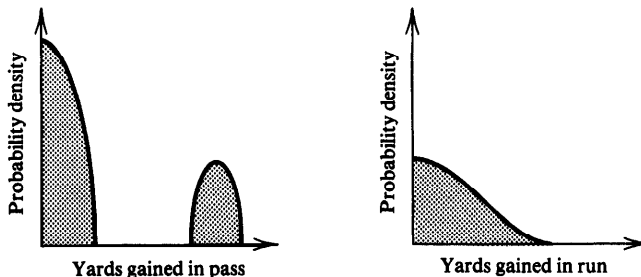


Figure 12.11. Distribution of yards gained in a run or a pass play.

To be precise, we wish to calculate the probability that  $\sum_{i=1}^n Z_i \geq n\alpha$ , where  $Z_i$  are independent random variables, and  $Z_i$  has a distribution corresponding to the strategy chosen.

The situation is illustrated in Figure 12.12. Let  $E$  be the set of types corresponding to the constraint,

$$E = \left\{ P: \sum_{a \in \mathcal{X}} P(a)a \geq \alpha \right\}. \quad (12.256)$$

If  $P_1$  is the distribution corresponding to passing all the time, then the probability of winning is the probability that the sample type is in  $E$ , which by Sanov's theorem is  $2^{-nD(P_1^* \| P_1)}$ , where  $P_1^*$  is the distribution in  $E$  that is closest to  $P_1$ . Similarly, if the coach uses the running game all the time, the probability of winning is  $2^{-nD(P_2^* \| P_2)}$ . What if he uses a mixture of strategies? Is it possible that  $2^{-nD(P_\lambda^* \| P_\lambda)}$ , the probability of winning with a mixed strategy,  $P_\lambda = \lambda P_1 + (1 - \lambda)P_2$ , is better than the probability of winning with either pure passing or pure running?

The somewhat surprising answer is yes, as can be shown by example. This provides a reason to use a mixed strategy other than the fact that it confuses the defense.

We end this section with another inequality due to Chernoff, which is a special version of Markov's inequality. This inequality is called the Chernoff bound.

**Lemma 12.9.1:** *Let  $Y$  be any random variable and let  $\psi(s)$  be the moment generating function of  $Y$ ,*

$$\psi(s) = Ee^{sY}. \quad (12.257)$$

*Then for all  $s \geq 0$ ,*

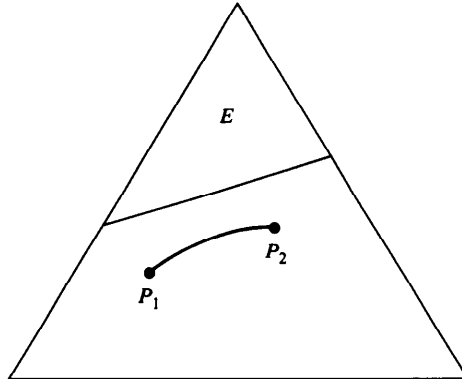


Figure 12.12. Probability simplex for a football game.



$$\Pr(Y \geq a) \leq e^{-sa} \psi(s). \quad (12.258)$$

**Proof:** Apply Markov's inequality to the non-negative random variable  $e^{sY}$ .  $\square$

## 12.10 LEMPEL-ZIV CODING

We now describe a scheme for universal data compression due to Ziv and Lempel [291], which is simple to implement and has an asymptotic rate approaching the entropy of the source. The algorithm is particularly simple and has become popular as the standard algorithm for file compression on computers because of its speed and efficiency.

We will consider a binary source throughout this section. The results generalize easily to any finite alphabet.

**Algorithm:** The source sequence is sequentially parsed into strings that have not appeared so far. For example, if the string is 1011010100010 . . . , we parse it as 1,0,11,01,010,00,10, . . . . After every comma, we look along the input sequence until we come to the shortest string that has not been marked off before. Since this is the shortest such string, all its prefixes must have occurred earlier. In particular, the string consisting of all but the last bit of this string must have occurred earlier. We code this phrase by giving the location of the prefix and the value of the last bit.

Let  $c(n)$  be the number of phrases in the parsing of the input  $n$ -sequence. We need  $\log c(n)$  bits to describe the location of the prefix to the phrase and 1 bit to describe the last bit. For example, the code for the above sequence is (000,1)(000,0)(001,1)(010,1)(100,0)(010,0)(001,0), where the first number of each pair gives the index of the prefix and the second number gives the last bit of the phrase. Decoding the coded sequence is straightforward and we can recover the source sequence without error.

The above algorithm requires two passes over the string—in the first pass, we parse the string and calculate  $c(n)$ , the number of phrases in the parsed string. We then use that to decide how many bits ( $\log c(n)$ ) to allot to the pointers in the algorithm. In the second pass, we calculate the pointers and produce the coded string as indicated above. The algorithm described above allots an equal number of bits to all the pointers. This is not necessary, since the range of the pointers is smaller at the initial portion of the string. The algorithm can be modified so that it requires only one pass over the string and uses fewer bits for the initial pointers. These modifications do not affect the asymptotic ef-

efficiency of the algorithm. Some of the implementation details are discussed by Welch [269] and Bell, Cleary and Witten [22].

In the example, we have not compressed the string; instead, we have expanded the number of bits by more than a factor of 2. But for long strings the phrases will get longer, and describing the phrases by describing the location of the prefix will be more efficient. We will show that this algorithm asymptotically achieves the entropy rate for the unknown ergodic source.

Without loss of generality, we will assume that the source alphabet is binary. Thus  $\mathcal{X} = \{0, 1\}$  throughout this section. We first define a parsing of the string to be a decomposition into phrases.

**Definition:** A parsing  $S$  of a binary string  $x_1 x_2 \dots x_n$  is a division of the string into phrases, separated by commas. A *distinct parsing* is a parsing such that no two phrases are identical.

For example, 0,111,1 is a distinct parsing of 01111, but 0,11,11 is a parsing which is not distinct.

The Lempel-Ziv algorithm described above gives a distinct parsing of the source sequence. Let  $c(n)$  denote the number of phrases in the Lempel-Ziv parsing of a sequence of length  $n$ . Of course,  $c(n)$  depends on the sequence  $X^n$ . The compressed sequence (after applying the Lempel-Ziv algorithm) consists of a list of  $c(n)$  pairs of numbers, each pair consisting of a pointer to the previous occurrence of the prefix of the phrase and the last bit of the phrase. Each pointer requires  $\log c(n)$  bits, and hence the total length of the compressed sequence is  $c(n)(\log c(n) + 1)$  bits. We will now show that  $\frac{c(n)(\log c(n) + 1)}{n} \rightarrow H(\mathcal{X})$  for a stationary ergodic sequence  $X_1, X_2, \dots, X_n$ . Our proof is based on the simple proof of asymptotic optimality of Lempel-Ziv coding due to Wyner and Ziv [285].

We first prove a few lemmas that we need for the proof of the theorem. The first is a bound on the number of phrases possible in a distinct parsing of a binary sequence of length  $n$ .

**Lemma 12.10.1 (Lempel and Ziv):** *The number of phrases  $c(n)$  in a distinct parsing of a binary sequence  $X_1, X_2, \dots, X_n$  satisfies*

$$c(n) \leq \frac{n}{(1 - \epsilon_n) \log n} \quad (12.259)$$

where  $\epsilon_n \rightarrow 0$  as  $n \rightarrow \infty$ .

**Proof:** Let

$$n_k = \sum_{j=1}^k j 2^j = (k-1)2^{k+1} + 2 \quad (12.260)$$

be the sum of the lengths of all distinct strings of length less than or equal to  $k$ . The number of phrases  $c$  in a distinct parsing of a sequence of length  $n$  is maximized when all the phrases are as short as possible. If  $n = n_k$ , this occurs when all the phrases are of length  $\leq k$ , and thus

$$c(n_k) \leq \sum_{j=1}^k 2^j = 2^{k+1} - 2 < 2^{k+1} \leq \frac{n_k}{k-1}. \quad (12.261)$$

If  $n_k \leq n < n_{k+1}$ , we write  $n = n_k + \Delta$ , where  $\Delta < (k+1)2^{k+1}$ . Then the parsing into shortest phrases has each of the phrases of length  $\leq k$  and  $\Delta/(k+1)$  phrases of length  $k+1$ . Thus

$$c(n) \leq \frac{n_k}{k-1} + \frac{\Delta}{k+1} \leq \frac{n_k + \Delta}{k-1} = \frac{n}{k-1}. \quad (12.262)$$

We now bound the size of  $k$  for a given  $n$ . Let  $n_k \leq n < n_{k+1}$ . Then

$$n \geq n_k = (k-1)2^{k+1} + 2 \geq 2^k, \quad (12.263)$$

and therefore

$$k \leq \log n. \quad (12.264)$$

Moreover,

$$n \leq n_{k+1} = k2^{k+2} + 2 \leq (k+2)2^{k+2} \leq (\log n + 2)2^{k+2} \quad (12.265)$$

by (12.264), and therefore

$$k+2 \geq \log \frac{n}{\log n + 2}, \quad (12.266)$$

or, for all  $n \geq 4$ ,

$$k-1 \geq \log n - \log(\log n + 2) - 3 \quad (12.267)$$

$$= \left(1 - \frac{\log(\log n + 2) + 3}{\log n}\right) \log n \quad (12.268)$$

$$\geq \left(1 - \frac{\log(2 \log n) + 3}{\log n}\right) \log n \quad (12.269)$$

$$= \left(1 - \frac{\log(\log n) + 4}{\log n}\right) \log n \quad (12.270)$$

$$= (1 - \epsilon_n) \log n. \quad (12.271)$$

Note that  $\epsilon_n = \min\{1, \frac{\log(\log n) + 4}{\log n}\}$ . Combining (12.271) with (12.262), we obtain the lemma.  $\square$

We will need a simple result on maximum entropy in the proof of the main theorem.

**Lemma 12.10.2:** *Let  $Z$  be a positive integer valued random variable with mean  $\mu$ . Then the entropy  $H(Z)$  is bounded by*

$$H(Z) \leq (\mu + 1) \log(\mu + 1) - \mu \log \mu. \quad (12.272)$$

**Proof:** The lemma follows directly from the results of Chapter 11, which show that the probability mass function that maximizes entropy subject to a constraint on the mean is the geometric distribution, for which we can compute the entropy. The details are left as an exercise for the reader.  $\square$

Let  $\{X_i\}_{i=-\infty}^{\infty}$  be a stationary ergodic process with probability mass function  $P(x_1, x_2, \dots, x_n)$ . (Ergodic processes are discussed in greater detail in Section 15.7.) For a fixed integer  $k$ , define the  $k$ th order Markov approximation to  $P$  as

$$Q_k(x_{-(k-1)}, \dots, x_0, x_1, \dots, x_n) \triangleq P(x_{-(k-1)}^0) \prod_{j=1}^n P(x_j | x_{j-k}^{j-1}), \quad (12.273)$$

where  $x_i^j \triangleq (x_i, x_{i+1}, \dots, x_j)$ ,  $i \leq j$ , and the initial state  $x_{-(k-1)}^0$  will be part of the specification of  $Q_k$ . Since  $P(X_n | X_{n-k}^{n-1})$  is itself an ergodic process, we have

$$-\frac{1}{n} \log Q_k(X_1, X_2, \dots, X_n | X_{-(k-1)}^0) = -\frac{1}{n} \sum_{j=1}^n \log P(X_j | X_{j-k}^{j-1}) \quad (12.274)$$

$$\rightarrow -E \log P(X_j | X_{j-k}^{j-1}) \quad (12.275)$$

$$= H(X_j | X_{j-k}^{j-1}). \quad (12.276)$$

We will bound the rate of the Lempel-Ziv code by the entropy rate of the  $k$ th order Markov approximation for all  $k$ . The entropy rate of the Markov approximation  $H(X_j | X_{j-k}^{j-1})$  converges to the entropy rate of the process as  $k \rightarrow \infty$  and this will prove the result.

Suppose  $X_{-(k-1)}^n = x_{-(k-1)}^n$ , and suppose that  $x_1^n$  is parsed into  $c$  distinct phrases,  $y_1, y_2, \dots, y_c$ . Let  $\nu_i$  be the index of the start of the  $i$ th phrase, i.e.,  $y_i = x_{\nu_i}^{\nu_i+1-1}$ . For each  $i = 1, 2, \dots, c$ , define  $s_i = x_{\nu_i-k}^{\nu_i-1}$ . Thus  $s_i$  is the  $k$  bits of  $x$  preceding  $y_i$ . Of course,  $s_1 = x_{-(k-1)}^0$ .

Let  $c_{ls}$  be the number of phrases  $y_i$  with length  $l$  and preceding state  $s_i = s$  for  $l = 1, 2, \dots$  and  $s \in \mathcal{X}^k$ . We then have

$$\sum_{l,s} c_{ls} = c \quad (12.277)$$

and

$$\sum_{l,s} lc_{ls} = n. \quad (12.278)$$

We now prove a surprising upper bound on the probability of a string based on the parsing of the string.

**Lemma 12.10.3** (*Ziv's inequality*): *For any distinct parsing (in particular, the Lempel-Ziv parsing) of the string  $x_1x_2\ldots x_n$ , we have*

$$\log Q_k(x_1, x_2, \ldots, x_n | s_1) \leq -\sum_{l,s} c_{ls} \log c_{ls}. \quad (12.279)$$

Note that the right hand side does not depend on  $Q_k$ .

**Proof:** We write

$$Q_k(x_1, x_2, \ldots, x_n | s_1) = Q(y_1, y_2, \ldots, y_c | s_1) \quad (12.280)$$

$$= \prod_{i=1}^c P(y_i | s_i), \quad (12.281)$$

or

$$\log Q_k(x_1, x_2, \ldots, x_n | s_1) = \sum_{i=1}^c \log P(y_i | s_i) \quad (12.282)$$

$$= \sum_{l,s} \sum_{i: |y_i|=l, s_i=s} \log P(y_i | s_i) \quad (12.283)$$

$$= \sum_{l,s} c_{ls} \sum_{i: |y_i|=l, s_i=s} \frac{1}{c_{ls}} \log P(y_i | s_i) \quad (12.284)$$

$$\leq \sum_{l,s} c_{ls} \log \left( \sum_{i: |y_i|=l, s_i=s} \frac{1}{c_{ls}} P(y_i | s_i) \right), \quad (12.285)$$

where the inequality follows from Jensen's inequality and the concavity of the logarithm.

Now since the  $y_i$  are distinct, we have  $\sum_{i: |y_i|=l, s_i=s} P(y_i | s_i) \leq 1$ . Thus

$$\log Q_k(x_1, x_2, \ldots, x_n | s_1) \leq \sum_{l,s} c_{ls} \log \frac{1}{c_{ls}}, \quad (12.286)$$

proving the lemma.  $\square$

We can now prove the main theorem:

**Theorem 12.10.1:** *Let  $\{X_n\}$  be a stationary ergodic process with entropy rate  $H(\mathcal{X})$ , and let  $c(n)$  be the number of phrases in a distinct parsing of a sample of length  $n$  from this process. Then*

$$\limsup_{n \rightarrow \infty} \frac{c(n) \log c(n)}{n} \leq H(\mathcal{X}) \quad (12.287)$$

with probability 1.

**Proof:** We will begin with Ziv's inequality, which we rewrite as

$$\log Q_k(x_1, x_2, \dots, x_n | s_1) \leq - \sum_{l,s} c_{ls} \log \frac{c_{ls} c}{c} \quad (12.288)$$

$$= -c \log c - c \sum_{ls} \frac{c_{ls}}{c} \log \frac{c_{ls}}{c}. \quad (12.289)$$

Writing  $\pi_{ls} = \frac{c_{ls}}{c}$ , we have

$$\sum_{l,s} \pi_{ls} = 1, \quad \sum_{l,s} l \pi_{ls} = \frac{n}{c}, \quad (12.290)$$

from (12.227) and (12.278). We now define random variables  $U, V$ , such that

$$\Pr(U = l, V = s) = \pi_{ls}. \quad (12.291)$$

Thus  $EU = \frac{n}{c}$  and

$$\log Q_k(x_1, x_2, \dots, x_n | s_1) \leq cH(U, V) - c \log c \quad (12.292)$$

or

$$-\frac{1}{n} \log Q_k(x_1, x_2, \dots, x_n | s_1) \geq \frac{c}{n} \log c - \frac{c}{n} H(U, V). \quad (12.293)$$

Now

$$H(U, V) \leq H(U) + H(V), \quad (12.294)$$

and  $H(V) \leq \log |\mathcal{X}|^k = k$ . By Lemma 12.10.2, we have

$$H(U) \leq (EU + 1) \log(EU + 1) - EU \log EU \quad (12.295)$$

$$= \left( \frac{n}{c} + 1 \right) \log \left( \frac{n}{c} + 1 \right) - \frac{n}{c} \log \frac{n}{c} \quad (12.296)$$

$$= \log \frac{n}{c} + \left( \frac{n}{c} + 1 \right) \log \left( \frac{c}{n} + 1 \right). \quad (12.297)$$

Thus

$$\frac{c}{n} H(U, V) \leq \frac{c}{n} k + \frac{c}{n} \log \frac{n}{c} + o(1). \quad (12.298)$$

For a given  $n$ , the maximum of  $\frac{c}{n} \log \frac{n}{c}$  is attained for the maximum value of  $c$  (for  $\frac{c}{n} \leq \frac{1}{e}$ ). But from Lemma 12.10.1,  $c \leq \frac{n}{\log n} (1 + o(1))$ . Thus

$$\frac{c}{n} \log \frac{n}{c} \leq O\left(\frac{\log \log n}{\log n}\right) \quad (12.299)$$

and therefore  $\frac{c}{n} H(U, V) \rightarrow 0$  as  $n \rightarrow \infty$ .

Therefore

$$\frac{c(n) \log c(n)}{n} \leq -\frac{1}{n} \log Q_k(x_1, x_2, \dots, x_n | s_1) + \epsilon_k(n) \quad (12.300)$$

where  $\epsilon_k(n) \rightarrow 0$  as  $n \rightarrow \infty$ . Hence, with probability 1,

$$\limsup_{n \rightarrow \infty} \frac{c(n) \log c(n)}{n} \leq \lim_{n \rightarrow \infty} -\frac{1}{n} \log Q_k(X_1, X_2, \dots, X_n | X_{-(k-1)}^0) \quad (12.301)$$

$$= H(X_0 | X_{-1}, \dots, X_{-k}) \quad (12.302)$$

$$\rightarrow H(\mathcal{X}) \quad \text{as } k \rightarrow \infty. \quad \square \quad (12.303)$$

We now prove that Lempel-Ziv coding is asymptotically optimal.

**Theorem 12.10.2:** *Let  $\{X_i\}_{-\infty}^{\infty}$  be a stationary ergodic stochastic process. Let  $l(X_1, X_2, \dots, X_n)$  be the Lempel-Ziv codeword length associated with  $X_1, X_2, \dots, X_n$ . Then*

$$\limsup_{n \rightarrow \infty} \frac{1}{n} l(X_1, X_2, \dots, X_n) \leq H(\mathcal{X}) \quad \text{with probability 1} \quad (12.304)$$

where  $H(\mathcal{X})$  is the entropy rate of the process.

**Proof:** We have shown that  $l(X_1, X_2, \dots, X_n) = c(n)(\log c(n) + 1)$ , where  $c(n)$  is the number of phrases in the Lempel-Ziv parsing of the string  $X_1, X_2, \dots, X_n$ . By Lemma 12.10.1,  $\limsup c(n)/n = 0$ , and thus Theorem 12.10.1 establishes that

$$\limsup \frac{l(X_1, X_2, \dots, X_n)}{n} = \limsup \left( \frac{c(n) \log c(n)}{n} + \frac{c(n)}{n} \right) \leq H(\mathcal{X}) \quad \text{with probability 1.} \quad \square \quad (12.305)$$

Thus the length per source symbol of the Lempel-Ziv encoding of an ergodic source is asymptotically no greater than the entropy rate of the source. The Lempel-Ziv code is a simple example of a universal code, i.e., a code that does not depend on the distribution of the source. This code can be used without knowledge of the source distribution and yet will achieve an asymptotic compression equal to the entropy rate of the source.

The Lempel-Ziv algorithm is now the standard algorithm for compression of files—it is implemented in the *compress* program in UNIX and in the *arc* program for PC's. The algorithm typically compresses ASCII text files by about a factor of 2. It has also been implemented in hardware and is used to effectively double the capacity of communication links for text files by compressing the file at one end and decompressing it at the other end.

## 12.11 FISHER INFORMATION AND THE CRAMÉR-RAO INEQUALITY

A standard problem in statistical estimation is to determine the parameters of a distribution from a sample of data drawn from that distribution. For example, let  $X_1, X_2, \dots, X_n$  be drawn i.i.d.  $\sim \mathcal{N}(\theta, 1)$ . Suppose we wish to estimate  $\theta$  from a sample of size  $n$ . There are a number of functions of the data that we can use to estimate  $\theta$ . For example, we can use the first sample  $X_1$ . Although the expected value of  $X_1$  is  $\theta$ , it is clear that we can do better by using more of the data. We guess that the best estimate of  $\theta$  is the sample mean  $\bar{X}_n = \frac{1}{n} \sum X_i$ . Indeed, it can be shown that  $\bar{X}_n$  is the minimum mean squared error unbiased estimator.

We begin with a few definitions. Let  $\{f(x; \theta)\}, \theta \in \Theta$ , denote an indexed family of densities,  $f(x; \theta) \geq 0$ ,  $\int f(x; \theta) dx = 1$  for all  $\theta \in \Theta$ . Here  $\Theta$  is called the *parameter set*.

**Definition:** An *estimator* for  $\theta$  for sample size  $n$  is a function  $T: \mathcal{X}^n \rightarrow \Theta$ .

An estimator is meant to approximate the value of the parameter. It is therefore desirable to have some idea of the goodness of the approximation. We will call the difference  $T - \theta$  the *error* of the estimator. The error is a random variable.

**Definition:** The *bias* of an estimator  $T(X_1, X_2, \dots, X_n)$  for the parameter  $\theta$  is the expected value of the error of the estimator, i.e., the bias is  $E_\theta T(X_1, X_2, \dots, X_n) - \theta$ . The subscript  $\theta$  means that the expectation is with respect to the density  $f(\cdot; \theta)$ . The estimator is said to be *unbiased* if



the bias is zero, i.e., the expected value of the estimator is equal to the parameter.

**Example 12.11.1:** Let  $X_1, X_2, \dots, X_n$  drawn i.i.d.  $\sim f(x) = (1/\lambda) e^{-x/\lambda}$ ,  $x \geq 0$  be a sequence of exponentially distributed random variables. Estimators of  $\lambda$  include  $X_1$  and  $\bar{X}_n$ . Both estimators are unbiased.

The bias is the expected value of the error, and the fact that it is zero does not guarantee that the error is low with high probability. We need to look at some loss function of the error; the most commonly chosen loss function is the expected square of the error. A good estimator should have a low expected squared error and should have an error that approaches 0 as the sample size goes to infinity. This motivates the following definition:

**Definition:** An estimator  $T(X_1, X_2, \dots, X_n)$  for  $\theta$  is said to be *consistent in probability* if  $T(X_1, X_2, \dots, X_n) \rightarrow \theta$  in probability as  $n \rightarrow \infty$ .

Consistency is a desirable asymptotic property, but we are interested in the behavior for small sample sizes as well. We can then rank estimators on the basis of their mean squared error.

**Definition:** An estimator  $T_1(X_1, X_2, \dots, X_n)$  is said to *dominate* another estimator  $T_2(X_1, X_2, \dots, X_n)$  if, for all  $\theta$ ,

$$E_\theta(T_1(X_1, X_2, \dots, X_n) - \theta)^2 \leq E_\theta(T_2(X_1, X_2, \dots, X_n) - \theta)^2, \quad (12.306)$$

This definition raises a natural question: what is the minimum variance unbiased estimator of  $\theta$ ? To answer this question, we derive the Cramér-Rao lower bound on the mean squared error of any estimator. We first define the score function of the distribution  $f(x; \theta)$ . We then use the Cauchy-Schwarz inequality to prove the Cramér-Rao lower bound on the variance of all unbiased estimators.

**Definition:** The *score*  $V$  is a random variable defined by

$$V = \frac{\partial}{\partial \theta} \ln f(X; \theta) = \frac{\frac{\partial}{\partial \theta} f(X; \theta)}{f(X; \theta)}, \quad (12.307)$$

where  $X \sim f(x; \theta)$ .

The mean value of the score is

$$EV = \int \frac{\frac{\partial}{\partial \theta} f(x; \theta)}{f(x; \theta)} f(x; \theta) dx \quad (12.308)$$

$$= \int \frac{\partial}{\partial \theta} f(x; \theta) dx \quad (12.309)$$

$$= \frac{\partial}{\partial \theta} \int f(x; \theta) dx \quad (12.310)$$

$$= \frac{\partial}{\partial \theta} 1 \quad (12.311)$$

$$= 0, \quad (12.312)$$

and therefore  $EV^2 = \text{var}(V)$ . The variance of the score has a special significance.

**Definition:** The Fisher information  $J(\theta)$  is the variance of the score, i.e.,

$$J(\theta) = E_{\theta} \left[ \frac{\partial}{\partial \theta} \ln f(x; \theta) \right]^2. \quad (12.313)$$

If we consider a sample of  $n$  random variables  $X_1, X_2, \dots, X_n$  drawn i.i.d.  $\sim f(x; \theta)$ , we have

$$f(x_1, x_2, \dots, x_n; \theta) = \prod_{i=1}^n f(x_i; \theta), \quad (12.314)$$

and the score function is the sum of the individual score functions,

$$V(X_1, X_2, \dots, X_n) = \frac{\partial}{\partial \theta} \ln f(X_1, X_2, \dots, X_n; \theta) \quad (12.315)$$

$$= \sum_{i=1}^n \frac{\partial}{\partial \theta} \ln f(X_i; \theta) \quad (12.316)$$

$$= \sum_{i=1}^n V(X_i), \quad (12.317)$$

where the  $V(X_i)$  are independent, identically distributed with zero mean. Hence the Fisher information is

$$J_n(\theta) = E_{\theta} \left[ \frac{\partial}{\partial \theta} \ln f(x_1, x_2, \dots, x_n; \theta) \right]^2 \quad (12.318)$$

$$= E_{\theta} V^2(X_1, X_2, \dots, X_n) \quad (12.319)$$

$$= E_{\theta} \left( \sum_{i=1}^n V(X_i) \right)^2 \quad (12.320)$$

$$= \sum_{i=1}^n E_{\theta} V^2(X_i) \quad (12.321)$$

$$= nJ(\theta). \quad (12.322)$$

Consequently, the Fisher information for  $n$  i.i.d. samples is  $n$  times the individual Fisher information. The significance of the Fisher information is shown in the following theorem:

**Theorem: 12.11.1** (*Cramér-Rao inequality*): *The mean squared error of any unbiased estimator  $T(X)$  of the parameter  $\theta$  is lower bounded by the reciprocal of the Fisher information, i.e.,*

$$\text{var}(T) \geq \frac{1}{J(\theta)}. \quad (12.323)$$

**Proof:** Let  $V$  be the score function and  $T$  be the estimator. By the Cauchy-Schwarz inequality, we have

$$(E_{\theta}[(V - E_{\theta}V)(T - E_{\theta}T)])^2 \leq E_{\theta}(V - E_{\theta}V)^2 E_{\theta}(T - E_{\theta}T)^2. \quad (12.324)$$

By (12.312),  $E_{\theta}V = 0$  and hence  $E_{\theta}(V - E_{\theta}V)(T - E_{\theta}T) = E_{\theta}(VT)$ . Also, by definition,  $\text{var}(V) = J(\theta)$ . Substituting these conditions in (12.324), we have

$$[E_{\theta}(VT)]^2 \leq J(\theta) \text{var}(T). \quad (12.325)$$

Now,

$$E_{\theta}(VT) = \int \frac{\frac{\partial}{\partial \theta} f(x; \theta)}{f(x; \theta)} T(x) f(x; \theta) dx \quad (12.326)$$

$$= \int \frac{\partial}{\partial \theta} f(x; \theta) T(x) dx \quad (12.327)$$

$$= \frac{\partial}{\partial \theta} \int f(x; \theta) T(x) dx \quad (12.328)$$

$$= \frac{\partial}{\partial \theta} E_{\theta}T \quad (12.329)$$

$$= \frac{\partial}{\partial \theta} \theta \quad (12.330)$$

$$= 1. \quad (12.331)$$

where the interchange of differentiation and integration in (12.328) can be justified using the bounded convergence theorem for appropriately well behaved  $f(x; \theta)$  and (12.330) follows from the fact that the estimator  $T$  is unbiased. Substituting this in (12.325), we obtain

$$\text{var}(T) \geq \frac{1}{J(\theta)}, \quad (12.332)$$

which is the Cramér-Rao inequality for unbiased estimators.  $\square$

By essentially the same arguments, we can show that for any estimator

$$E_{\theta}(T - \theta)^2 \geq \frac{[1 + b'_T(\theta)]^2}{J(\theta)} + b_T^2(\theta), \quad (12.333)$$

where  $b_T(\theta) = E_{\theta}T - \theta$  and  $b'_T(\theta)$  is the derivative of  $b_T(\theta)$  with respect to  $\theta$ . The proof of this is left as an exercise at the end of the chapter.

**Example 12.11.2:** Let  $X_1, X_2, \dots, X_n$  be i.i.d.  $\sim \mathcal{N}(\theta, \sigma^2)$ ,  $\sigma^2$  known. Here  $J(\theta) = \frac{n}{\sigma^2}$ . Let  $T(X_1, X_2, \dots, X_n) = \bar{X}_n = \frac{1}{n} \sum X_i$ . Then  $E_{\theta}(\bar{X}_n - \theta)^2 = \frac{\sigma^2}{n} = \frac{1}{J(\theta)}$ . Thus  $\bar{X}_n$  is the minimum variance unbiased estimator of  $\theta$ , since it achieves the Cramér-Rao lower bound.

The Cramér-Rao inequality gives us the lowest possible variance for all unbiased estimators. We now use it to define the most efficient estimator.

**Definition:** An unbiased estimator  $T$  is said to be *efficient* if it meets the Cramér-Rao bound with equality, i.e., if  $\text{var}(T) = \frac{1}{J(\theta)}$ .

The Fisher information is therefore a measure of the amount of “information” about  $\theta$  that is present in the data. It gives a lower bound on the error in estimating  $\theta$  from the data. However, it is possible that there does not exist an estimator meeting this lower bound.

We can generalize the concept of Fisher information to the multi-parameter case, in which case we define the Fisher information matrix  $J(\theta)$  with elements

$$J_{ij}(\theta) = \int f(x; \theta) \frac{\partial}{\partial \theta_i} \ln f(x; \theta) \frac{\partial}{\partial \theta_j} \ln f(x; \theta) dx. \quad (12.334)$$

The Cramér-Rao inequality becomes the matrix inequality

$$\Sigma \geq J^{-1}(\theta), \quad (12.335)$$

where  $\Sigma$  is the covariance matrix of a set of unbiased estimators for the parameters  $\theta$  and  $\Sigma \geq J^{-1}(\theta)$  in the sense that the difference  $\Sigma - J^{-1}$  is a non-negative definite matrix. We will not go into the details of the proof for multiple parameters; the basic ideas are similar.

Is there a relationship between the Fisher information  $J(\theta)$  and quantities like entropy defined earlier? Note that Fisher information is defined with respect to a family of parametric distributions, unlike entropy, which is defined for all distributions. But we can parametrize any distribution,  $f(x)$ , by a location parameter  $\theta$  and define Fisher information with respect to the family of densities  $\{f(x - \theta)\}$  under translation. We will explore the relationship in greater detail in Section 16.7, where we show that while entropy is related to the volume of the typical set, the Fisher information is related to the surface area of the typical set. Further relationships of Fisher information to relative entropy are developed in the exercises.

### SUMMARY OF CHAPTER 12

#### Basic identities:

$$Q^n(\mathbf{x}) = 2^{-n(D(P_{\mathbf{x}}\|Q) + H(P_{\mathbf{x}}))}, \quad (12.336)$$

$$|\mathcal{P}_n| \leq (n+1)^{|\mathcal{X}|}, \quad (12.337)$$

$$|T(P)| \doteq 2^{nH(P)}, \quad (12.338)$$

$$Q^n(T(P)) \doteq 2^{-nD(P\|Q)}. \quad (12.339)$$

#### Universal data compression:

$$P_e^{(n)} \leq 2^{-nD(P_R^*\|Q)}, \quad \text{for all } Q, \quad (12.340)$$

where

$$D(P_R^*\|Q) = \min_{P: H(P) \geq R} D(P\|Q) \quad (12.341)$$

#### Large deviations (Sanov's theorem):

$$Q^n(E) = Q^n(E \cap \mathcal{P}_n) \leq (n+1)^{|\mathcal{X}|} 2^{-nD(P^*\|Q)}, \quad (12.342)$$

$$D(P^*\|Q) = \min_{P \in E} D(P\|Q), \quad (12.343)$$

If  $E$  is the closure of its interior, then

$$Q^n(E) \doteq 2^{-nD(P^*\|Q)}. \quad (12.344)$$

#### $\mathcal{L}_1$ bound on relative entropy:

$$D(P_1\|P_2) \geq \frac{1}{2 \ln 2} \|P_1 - P_2\|_1^2 \quad (12.345)$$

**Pythagorean theorem:** If  $E$  is a convex set of types, distribution  $Q \notin E$ , and  $P^*$  achieves  $D(P^*||Q) = \min_{P \in E} D(P||Q)$ , we have

$$D(P||Q) \geq D(P||P^*) + D(P^*||Q) \quad (12.346)$$

for all  $P \in E$ .

**Conditional limit theorem:** If  $X_1, X_2, \dots, X_n$  i.i.d.  $\sim Q$ , then

$$\Pr(X_1 = a | P_{X^n} \in E) \rightarrow P^*(a) \text{ in probability,} \quad (12.347)$$

where  $P^*$  minimizes  $D(P||Q)$  over  $P \in E$ . In particular,

$$\Pr\left\{X_1 = a \mid \frac{1}{n} \sum_{i=1}^n X_i \geq \alpha\right\} \rightarrow \frac{Q(a)e^{\lambda \alpha}}{\sum_x Q(x)e^{\lambda x}}. \quad (12.348)$$

**Neyman-Pearson lemma:** The optimum test between two densities  $P_1$  and  $P_2$  has a decision region of the form "Accept  $P = P_1$  if  $\frac{P_1(x_1, x_2, \dots, x_n)}{P_2(x_1, x_2, \dots, x_n)} > T$ ."

**Stein's lemma:** The best achievable error exponent  $\beta_n^*$  if  $\alpha_n \leq \epsilon$ :

$$\beta_n^* = \min_{\substack{A_n \subseteq \mathcal{X}^n \\ \alpha_n < \epsilon}} \beta_n. \quad (12.349)$$

$$\lim_{n \rightarrow \infty} \lim_{\epsilon \rightarrow 0} \frac{1}{n} \log \beta_n^* = -D(P_1||P_2). \quad (12.350)$$

**Chernoff information:** The best achievable exponent for a Bayesian probability of error is

$$D^* = D(P_{\lambda^*}||P_1) = D(P_{\lambda^*}||P_2), \quad (12.351)$$

where

$$P_{\lambda} = \frac{P_1^{\lambda}(x)P_2^{1-\lambda}(x)}{\sum_{a \in \mathcal{X}} P_1^{\lambda}(a)P_2^{1-\lambda}(a)} \quad (12.352)$$

with  $\lambda = \lambda^*$  chosen so that

$$D(P_{\lambda}||P_1) = D(P_{\lambda}||P_2). \quad (12.353)$$

**Lempel-Ziv:** Universal data compression. For a stationary ergodic source,

$$\limsup \frac{l(X_1, X_2, \dots, X_n)}{n} = \limsup \frac{c(n) \log c(n)}{n} \leq H(\mathcal{X}). \quad (12.354)$$

**Fisher information:**

$$J(\theta) = E_{\theta} \left[ \frac{\partial}{\partial \theta} \ln f(x; \theta) \right]^2. \quad (12.355)$$

**Cramér-Rao inequality:** For any unbiased estimator  $T$  of  $\theta$ ,

$$E_{\theta}(T(X) - \theta)^2 = \text{var}(T) \geq \frac{1}{J(\theta)}. \quad (12.356)$$

## PROBLEMS FOR CHAPTER 12

1. *Stein's lemma.* Consider the two hypothesis test

$$H_1: f = f_1 \quad \text{vs.} \quad H_2: f = f_2$$

Find  $D(f_1 \| f_2)$  if

- $f_i(x) = N(0, \sigma_i^2)$ ,  $i = 1, 2$
  - $f_i(x) = \lambda_i e^{-\lambda_i x}$ ,  $x \geq 0$ ,  $i = 1, 2$
  - $f_1(x)$  is the uniform density over the interval  $[0, 1]$  and  $f_2(x)$  is the uniform density over  $[a, a + 1]$ . Assume  $0 < a < 1$ .
  - $f_1$  corresponds to a fair coin and  $f_2$  corresponds to a two-headed coin.
2. A relation between  $D(P \| Q)$  and chi-square. Show that the  $\chi^2$  statistic

$$\chi^2 = \sum_x \frac{(P(x) - Q(x))^2}{Q(x)}$$

is (twice) the first term in the Taylor series expansion of  $D(P \| Q)$  about  $Q$ . Thus  $D(P \| Q) = \frac{1}{2} \chi^2 + \dots$ .

*Hint:* Write  $\frac{P}{Q} = 1 + \frac{P-Q}{Q}$  and expand the log.

3. *Error exponent for universal codes.* A universal source code of rate  $R$  achieves a probability of error  $P_e^{(n)} \doteq e^{-nD(P^* \| Q)}$ , where  $Q$  is the true distribution and  $P^*$  achieves  $\min D(P \| Q)$  over all  $P$  such that  $H(P) \geq R$ .
- Find  $P^*$  in terms of  $Q$  and  $R$ .
  - Now let  $X$  be binary. Find the region of source probabilities  $Q(x)$ ,  $x \in \{0, 1\}$ , for which rate  $R$  is sufficient for the universal source code to achieve  $P_e^{(n)} \rightarrow 0$ .
4. *Sequential projection.* We wish to show that projecting  $Q$  onto  $P_1$  and then projecting the projection  $\hat{Q}$  onto  $P_1 \cap P_2$  is the same as projecting  $Q$  directly onto  $P_1 \cap P_2$ . Let  $\mathcal{P}_1$  be the set of probability mass functions on  $\mathcal{X}$  satisfying

$$\sum_x p(x) = 1, \quad (12.357)$$

$$\sum_x p(x) h_i(x) \geq \alpha_i, \quad i = 1, 2, \dots, r. \quad (12.358)$$

Let  $\mathcal{P}_2$  be the set of probability mass functions on  $\mathcal{X}$  satisfying

$$\sum_x p(x) = 1, \quad (12.359)$$

$$\sum_x p(x) g_j(x) \geq \beta_j, \quad j = 1, 2, \dots, s. \quad (12.360)$$

Suppose  $Q \notin P_1 \cup P_2$ . Let  $P^*$  minimize  $D(P \| Q)$  over all  $P \in \mathcal{P}_1$ . Let  $R^*$  minimize  $D(R \| Q)$  over all  $R \in \mathcal{P}_1 \cap \mathcal{P}_2$ . Argue that  $R^*$  minimizes  $D(R \| P^*)$  over all  $R \in P_1 \cap P_2$ .

5. *Counting.* Let  $\mathcal{X} = \{1, 2, \dots, m\}$ . Show that the number of sequences  $x^n \in \mathcal{X}^n$  satisfying  $\frac{1}{n} \sum_{i=1}^n g(x_i) \geq \alpha$  is approximately equal to  $2^{nH^*}$ , to first order in the exponent, for  $n$  sufficiently large, where

$$H^* = \max_{P: \sum_{i=1}^m P(i)g(i) \geq \alpha} H(P). \quad (12.361)$$

6. *Biased estimates may be better.* Consider the problem of estimating  $\mu$  and  $\sigma^2$  from  $n$  samples of data drawn i.i.d. from a  $\mathcal{N}(\mu, \sigma^2)$  distribution.

- (a) Show that  $\bar{X}$  is an unbiased estimator of  $\mu$ .  
 (b) Show that the estimator

$$S_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \quad (12.362)$$

is biased and the estimator

$$S_{n-1}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \quad (12.363)$$

is unbiased.

- (c) Show that  $S_n^2$  has a lower mean squared error than  $S_{n-1}^2$ . This illustrates the idea that a biased estimator may be “better” than an unbiased estimator.
7. *Fisher information and relative entropy.* Show for a parametric family  $\{p_\theta(x)\}$  that

$$\lim_{\theta' \rightarrow \theta} \frac{1}{(\theta - \theta')^2} D(p_\theta \| p_{\theta'}) = \frac{1}{\ln 4} J(\theta). \quad (12.364)$$

8. *Examples of Fisher information.* The Fisher information  $J(\theta)$  for the family  $f_\theta(x)$ ,  $\theta \in \mathbb{R}$  is defined by

$$J(\theta) = E_\theta \left( \frac{\partial f_\theta(X)/\partial \theta}{f_\theta(X)} \right)^2 = \int \frac{(f'_\theta)^2}{f_\theta}.$$

Find the Fisher information for the following families:

- (a)  $f_\theta(x) = N(0, \theta) = \frac{1}{\sqrt{2\pi\theta}} e^{-\frac{x^2}{2\theta}}$   
 (b)  $f_\theta(x) = \theta e^{-\theta x}$ ,  $x \geq 0$   
 (c) What is the Cramér-Rao lower bound on  $E_\theta(\hat{\theta}(X) - \theta)^2$ , where  $\hat{\theta}(X)$  is an unbiased estimator of  $\theta$  for (a) and (b)?
9. *Two conditionally independent looks double the Fisher information.* Let  $g_\theta(x_1, x_2) = f_\theta(x_1)f_\theta(x_2)$ . Show  $J_g(\theta) = 2J_f(\theta)$ .
10. *Joint distributions and product distributions.* Consider a joint distribution  $Q(x, y)$  with marginals  $Q(x)$  and  $Q(y)$ . Let  $E$  be the set of types that look jointly typical with respect to  $Q$ , i.e.,



$$\begin{aligned}
E = \{ & P(x, y) : - \sum_{x, y} P(x, y) \log Q(x) - H(X) = 0, \\
& - \sum_{x, y} P(x, y) \log Q(y) - H(Y) = 0, \\
& - \sum_{x, y} P(x, y) \log Q(x, y) - H(X, Y) = 0 \} . \quad (12.365)
\end{aligned}$$

- (a) Let  $Q_0(x, y)$  be another distribution on  $\mathcal{X} \times \mathcal{Y}$ . Argue that the distribution  $P^*$  in  $E$  that is closest to  $Q_0$  is of the form

$$P^*(x, y) = Q_0(x, y) e^{\lambda_0 + \lambda_1 \log Q(x) + \lambda_2 \log Q(y) + \lambda_3 \log Q(x, y)}, \quad (12.366)$$

where  $\lambda_0, \lambda_1, \lambda_2$  and  $\lambda_3$  are chosen to satisfy the constraints. Argue that this distribution is unique.

- (b) Now let  $Q_0(x, y) = Q(x)Q(y)$ . Verify that  $Q(x, y)$  is of the form (12.366) and satisfies the constraints. Thus  $P^*(x, y) = Q(x, y)$ , i.e., the distribution in  $E$  closest to the product distribution is the joint distribution.
11. *Cramér-Rao inequality with a bias term.* Let  $X \sim f(x; \theta)$  and let  $T(X)$  be an estimator for  $\theta$ . Let  $b_T(\theta) = E_\theta T - \theta$  be the bias of the estimator. Show that

$$E_\theta (T - \theta)^2 \geq \frac{[1 + b_T'(\theta)]^2}{J(\theta)} + b_T^2(\theta). \quad (12.367)$$

12. *Lempel-Ziv.* Give the Lempel-Ziv parsing and encoding of 000000110-10100000110101.

## HISTORICAL NOTES

The method of types evolved from notions of weak typicality and strong typicality; some of the ideas were used by Wolfowitz [277] to prove channel capacity theorems. The method was fully developed by Csiszár and Körner [83], who derived the main theorems of information theory from this viewpoint. The method of types described in Section 12.1 follows the development in Csiszár and Körner. The  $\mathcal{L}_1$  lower bound on relative entropy is due to Csiszár [78], Kullback [151] and Kemperman [227]. Sanov's theorem [175] was generalized by Csiszár [289] using the method of types.

The parsing algorithm for Lempel-Ziv encoding was introduced by Lempel and Ziv [175] and was proved to achieve the entropy rate by Ziv [289]. The algorithm described in the text was first described in Ziv and Lempel [289]. A more transparent proof was provided by Wyner and Ziv [285], which we have used to prove the results in Section 12.10. A number of different variations of the basic Lempel-Ziv algorithm are described in the book by Bell, Cleary and Witten [22].