

- (b) First solve the inference problem of determining the conditional density  $p(t|\mathbf{x})$ , and then subsequently marginalize to find the conditional mean given by (1.89).
- (c) Find a regression function  $y(\mathbf{x})$  directly from the training data.

The relative merits of these three approaches follow the same lines as for classification problems above.

The squared loss is not the only possible choice of loss function for regression. Indeed, there are situations in which squared loss can lead to very poor results and where we need to develop more sophisticated approaches. An important example concerns situations in which the conditional distribution  $p(t|\mathbf{x})$  is multimodal, as often arises in the solution of inverse problems. Here we consider briefly one simple generalization of the squared loss, called the *Minkowski* loss, whose expectation is given by

$$\mathbb{E}[L_q] = \iint |y(\mathbf{x}) - t|^q p(\mathbf{x}, t) d\mathbf{x} dt \quad (1.91)$$

which reduces to the expected squared loss for  $q = 2$ . The function  $|y - t|^q$  is plotted against  $y - t$  for various values of  $q$  in Figure 1.29. The minimum of  $\mathbb{E}[L_q]$  is given by the conditional mean for  $q = 2$ , the conditional median for  $q = 1$ , and the conditional mode for  $q \rightarrow 0$ .

Section 5.6

Exercise 1.27

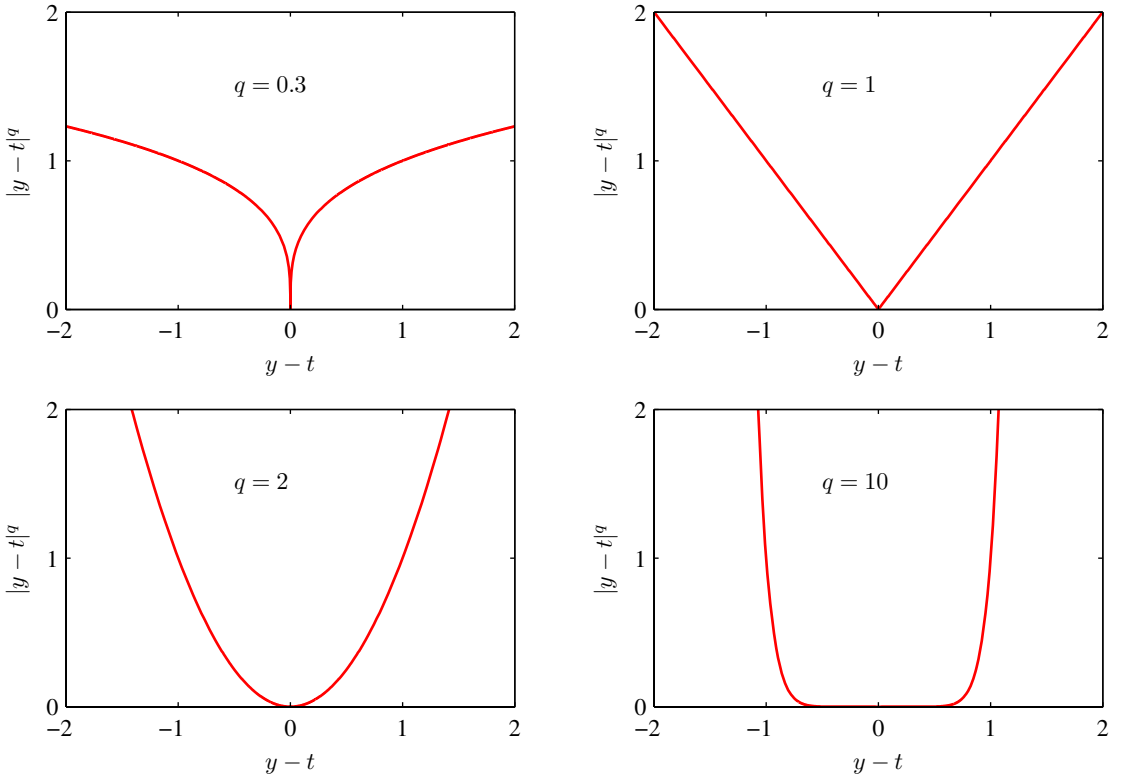
## 1.6. Information Theory

---

In this chapter, we have discussed a variety of concepts from probability theory and decision theory that will form the foundations for much of the subsequent discussion in this book. We close this chapter by introducing some additional concepts from the field of information theory, which will also prove useful in our development of pattern recognition and machine learning techniques. Again, we shall focus only on the key concepts, and we refer the reader elsewhere for more detailed discussions (Viterbi and Omura, 1979; Cover and Thomas, 1991; MacKay, 2003).

We begin by considering a discrete random variable  $x$  and we ask how much information is received when we observe a specific value for this variable. The amount of information can be viewed as the ‘degree of surprise’ on learning the value of  $x$ . If we are told that a highly improbable event has just occurred, we will have received more information than if we were told that some very likely event has just occurred, and if we knew that the event was certain to happen we would receive no information. Our measure of information content will therefore depend on the probability distribution  $p(x)$ , and we therefore look for a quantity  $h(x)$  that is a monotonic function of the probability  $p(x)$  and that expresses the information content. The form of  $h(\cdot)$  can be found by noting that if we have two events  $x$  and  $y$  that are unrelated, then the information gain from observing both of them should be the sum of the information gained from each of them separately, so that  $h(x, y) = h(x) + h(y)$ . Two unrelated events will be statistically independent and so  $p(x, y) = p(x)p(y)$ . From these two relationships, it is easily shown that  $h(x)$  must be given by the logarithm of  $p(x)$  and so we have

Exercise 1.28



**Figure 1.29** Plots of the quantity  $L_q = |y - t|^q$  for various values of  $q$ .

$$h(x) = -\log_2 p(x) \quad (1.92)$$

where the negative sign ensures that information is positive or zero. Note that low probability events  $x$  correspond to high information content. The choice of basis for the logarithm is arbitrary, and for the moment we shall adopt the convention prevalent in information theory of using logarithms to the base of 2. In this case, as we shall see shortly, the units of  $h(x)$  are bits ('binary digits').

Now suppose that a sender wishes to transmit the value of a random variable to a receiver. The average amount of information that they transmit in the process is obtained by taking the expectation of (1.92) with respect to the distribution  $p(x)$  and is given by

$$H[x] = -\sum_x p(x) \log_2 p(x). \quad (1.93)$$

This important quantity is called the *entropy* of the random variable  $x$ . Note that  $\lim_{p \rightarrow 0} p \ln p = 0$  and so we shall take  $p(x) \ln p(x) = 0$  whenever we encounter a value for  $x$  such that  $p(x) = 0$ .

So far we have given a rather heuristic motivation for the definition of informa-

tion (1.92) and the corresponding entropy (1.93). We now show that these definitions indeed possess useful properties. Consider a random variable  $x$  having 8 possible states, each of which is equally likely. In order to communicate the value of  $x$  to a receiver, we would need to transmit a message of length 3 bits. Notice that the entropy of this variable is given by

$$H[x] = -8 \times \frac{1}{8} \log_2 \frac{1}{8} = 3 \text{ bits.}$$

Now consider an example (Cover and Thomas, 1991) of a variable having 8 possible states  $\{a, b, c, d, e, f, g, h\}$  for which the respective probabilities are given by  $(\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64})$ . The entropy in this case is given by

$$H[x] = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{16} \log_2 \frac{1}{16} - \frac{4}{64} \log_2 \frac{1}{64} = 2 \text{ bits.}$$

We see that the nonuniform distribution has a smaller entropy than the uniform one, and we shall gain some insight into this shortly when we discuss the interpretation of entropy in terms of disorder. For the moment, let us consider how we would transmit the identity of the variable's state to a receiver. We could do this, as before, using a 3-bit number. However, we can take advantage of the nonuniform distribution by using shorter codes for the more probable events, at the expense of longer codes for the less probable events, in the hope of getting a shorter average code length. This can be done by representing the states  $\{a, b, c, d, e, f, g, h\}$  using, for instance, the following set of code strings: 0, 10, 110, 1110, 111100, 111101, 111110, 111111. The average length of the code that has to be transmitted is then

$$\text{average code length} = \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{16} \times 4 + 4 \times \frac{1}{64} \times 6 = 2 \text{ bits}$$

which again is the same as the entropy of the random variable. Note that shorter code strings cannot be used because it must be possible to disambiguate a concatenation of such strings into its component parts. For instance, 11001110 decodes uniquely into the state sequence  $c, a, d$ .

This relation between entropy and shortest coding length is a general one. The *noiseless coding theorem* (Shannon, 1948) states that the entropy is a lower bound on the number of bits needed to transmit the state of a random variable.

From now on, we shall switch to the use of natural logarithms in defining entropy, as this will provide a more convenient link with ideas elsewhere in this book. In this case, the entropy is measured in units of 'nats' instead of bits, which differ simply by a factor of  $\ln 2$ .

We have introduced the concept of entropy in terms of the average amount of information needed to specify the state of a random variable. In fact, the concept of entropy has much earlier origins in physics where it was introduced in the context of equilibrium thermodynamics and later given a deeper interpretation as a measure of disorder through developments in statistical mechanics. We can understand this alternative view of entropy by considering a set of  $N$  identical objects that are to be divided amongst a set of bins, such that there are  $n_i$  objects in the  $i^{\text{th}}$  bin. Consider

the number of different ways of allocating the objects to the bins. There are  $N$  ways to choose the first object,  $(N - 1)$  ways to choose the second object, and so on, leading to a total of  $N!$  ways to allocate all  $N$  objects to the bins, where  $N!$  (pronounced ‘factorial  $N$ ’) denotes the product  $N \times (N - 1) \times \cdots \times 2 \times 1$ . However, we don’t wish to distinguish between rearrangements of objects within each bin. In the  $i^{\text{th}}$  bin there are  $n_i!$  ways of reordering the objects, and so the total number of ways of allocating the  $N$  objects to the bins is given by

$$W = \frac{N!}{\prod_i n_i!} \quad (1.94)$$

which is called the *multiplicity*. The entropy is then defined as the logarithm of the multiplicity scaled by an appropriate constant

$$H = \frac{1}{N} \ln W = \frac{1}{N} \ln N! - \frac{1}{N} \sum_i \ln n_i! \quad (1.95)$$

We now consider the limit  $N \rightarrow \infty$ , in which the fractions  $n_i/N$  are held fixed, and apply Stirling’s approximation

$$\ln N! \simeq N \ln N - N \quad (1.96)$$

which gives

$$H = - \lim_{N \rightarrow \infty} \sum_i \left( \frac{n_i}{N} \right) \ln \left( \frac{n_i}{N} \right) = - \sum_i p_i \ln p_i \quad (1.97)$$

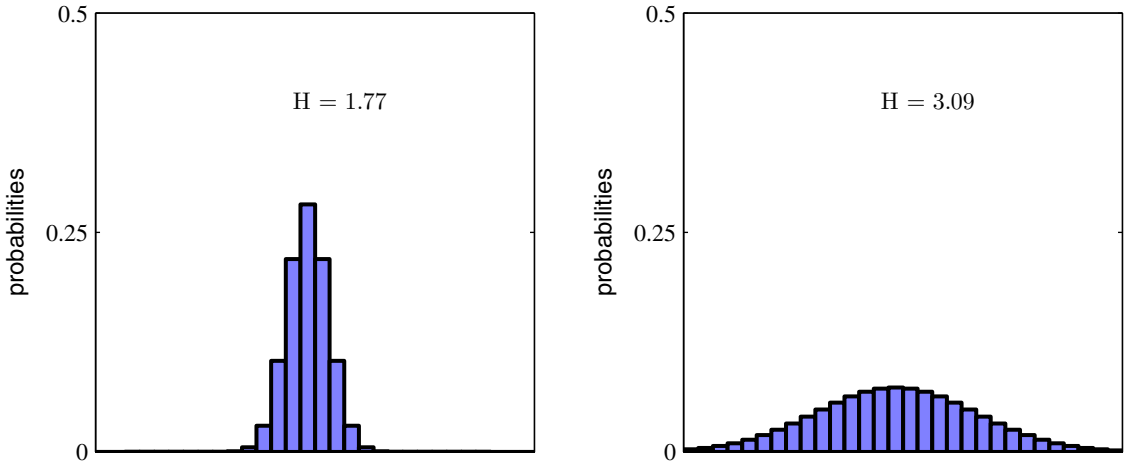
where we have used  $\sum_i n_i = N$ . Here  $p_i = \lim_{N \rightarrow \infty} (n_i/N)$  is the probability of an object being assigned to the  $i^{\text{th}}$  bin. In physics terminology, the specific arrangements of objects in the bins is called a *microstate*, and the overall distribution of occupation numbers, expressed through the ratios  $n_i/N$ , is called a *macrostate*. The multiplicity  $W$  is also known as the *weight* of the macrostate.

We can interpret the bins as the states  $x_i$  of a discrete random variable  $X$ , where  $p(X = x_i) = p_i$ . The entropy of the random variable  $X$  is then

$$H[p] = - \sum_i p(x_i) \ln p(x_i). \quad (1.98)$$

Distributions  $p(x_i)$  that are sharply peaked around a few values will have a relatively low entropy, whereas those that are spread more evenly across many values will have higher entropy, as illustrated in Figure 1.30. Because  $0 \leq p_i \leq 1$ , the entropy is nonnegative, and it will equal its minimum value of 0 when one of the  $p_i = 1$  and all other  $p_{j \neq i} = 0$ . The maximum entropy configuration can be found by maximizing  $H$  using a Lagrange multiplier to enforce the normalization constraint on the probabilities. Thus we maximize

$$\tilde{H} = - \sum_i p(x_i) \ln p(x_i) + \lambda \left( \sum_i p(x_i) - 1 \right) \quad (1.99)$$



**Figure 1.30** Histograms of two probability distributions over 30 bins illustrating the higher value of the entropy  $H$  for the broader distribution. The largest entropy would arise from a uniform distribution that would give  $H = -\ln(1/30) = 3.40$ .

from which we find that all of the  $p(x_i)$  are equal and are given by  $p(x_i) = 1/M$  where  $M$  is the total number of states  $x_i$ . The corresponding value of the entropy is then  $H = \ln M$ . This result can also be derived from Jensen's inequality (to be discussed shortly). To verify that the stationary point is indeed a maximum, we can evaluate the second derivative of the entropy, which gives

*Exercise 1.29*

$$\frac{\partial \tilde{H}}{\partial p(x_i) \partial p(x_j)} = -I_{ij} \frac{1}{p_i} \quad (1.100)$$

where  $I_{ij}$  are the elements of the identity matrix.

We can extend the definition of entropy to include distributions  $p(x)$  over continuous variables  $x$  as follows. First divide  $x$  into bins of width  $\Delta$ . Then, assuming  $p(x)$  is continuous, the *mean value theorem* (Weisstein, 1999) tells us that, for each such bin, there must exist a value  $x_i$  such that

$$\int_{i\Delta}^{(i+1)\Delta} p(x) dx = p(x_i) \Delta. \quad (1.101)$$

We can now quantize the continuous variable  $x$  by assigning any value  $x$  to the value  $x_i$  whenever  $x$  falls in the  $i^{\text{th}}$  bin. The probability of observing the value  $x_i$  is then  $p(x_i) \Delta$ . This gives a discrete distribution for which the entropy takes the form

$$H_{\Delta} = - \sum_i p(x_i) \Delta \ln(p(x_i) \Delta) = - \sum_i p(x_i) \Delta \ln p(x_i) - \ln \Delta \quad (1.102)$$

where we have used  $\sum_i p(x_i) \Delta = 1$ , which follows from (1.101). We now omit the second term  $-\ln \Delta$  on the right-hand side of (1.102) and then consider the limit

$\Delta \rightarrow 0$ . The first term on the right-hand side of (1.102) will approach the integral of  $p(x) \ln p(x)$  in this limit so that

$$\lim_{\Delta \rightarrow 0} \left\{ \sum_i p(x_i) \Delta \ln p(x_i) \right\} = - \int p(x) \ln p(x) dx \quad (1.103)$$

where the quantity on the right-hand side is called the *differential entropy*. We see that the discrete and continuous forms of the entropy differ by a quantity  $\ln \Delta$ , which diverges in the limit  $\Delta \rightarrow 0$ . This reflects the fact that to specify a continuous variable very precisely requires a large number of bits. For a density defined over multiple continuous variables, denoted collectively by the vector  $\mathbf{x}$ , the differential entropy is given by

$$H[\mathbf{x}] = - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x}. \quad (1.104)$$

In the case of discrete distributions, we saw that the maximum entropy configuration corresponded to an equal distribution of probabilities across the possible states of the variable. Let us now consider the maximum entropy configuration for a continuous variable. In order for this maximum to be well defined, it will be necessary to constrain the first and second moments of  $p(x)$  as well as preserving the normalization constraint. We therefore maximize the differential entropy with the



**Ludwig Boltzmann**  
1844–1906

Ludwig Eduard Boltzmann was an Austrian physicist who created the field of statistical mechanics. Prior to Boltzmann, the concept of entropy was already known from classical thermodynamics where it quantifies the fact that when we take energy from a system, not all of that energy is typically available to do useful work. Boltzmann showed that the thermodynamic entropy  $S$ , a macroscopic quantity, could be related to the statistical properties at the microscopic level. This is expressed through the famous equation  $S = k \ln W$  in which  $W$  represents the number of possible microstates in a macrostate, and  $k \simeq 1.38 \times 10^{-23}$  (in units of Joules per Kelvin) is known as Boltzmann's constant. Boltzmann's ideas were disputed by many scientists of their day. One difficulty they saw arose from the second law of thermo-

dynamics, which states that the entropy of a closed system tends to increase with time. By contrast, at the microscopic level the classical Newtonian equations of physics are reversible, and so they found it difficult to see how the latter could explain the former. They didn't fully appreciate Boltzmann's arguments, which were statistical in nature and which concluded not that entropy could never decrease over time but simply that with overwhelming probability it would generally increase. Boltzmann even had a long-running dispute with the editor of the leading German physics journal who refused to let him refer to atoms and molecules as anything other than convenient theoretical constructs. The continued attacks on his work led to bouts of depression, and eventually he committed suicide. Shortly after Boltzmann's death, new experiments by Perrin on colloidal suspensions verified his theories and confirmed the value of the Boltzmann constant. The equation  $S = k \ln W$  is carved on Boltzmann's tombstone.

three constraints

$$\int_{-\infty}^{\infty} p(x) dx = 1 \quad (1.105)$$

$$\int_{-\infty}^{\infty} xp(x) dx = \mu \quad (1.106)$$

$$\int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx = \sigma^2. \quad (1.107)$$

### Appendix E

The constrained maximization can be performed using Lagrange multipliers so that we maximize the following functional with respect to  $p(x)$

$$\begin{aligned} & - \int_{-\infty}^{\infty} p(x) \ln p(x) dx + \lambda_1 \left( \int_{-\infty}^{\infty} p(x) dx - 1 \right) \\ & + \lambda_2 \left( \int_{-\infty}^{\infty} xp(x) dx - \mu \right) + \lambda_3 \left( \int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx - \sigma^2 \right). \end{aligned}$$

### Appendix D

Using the calculus of variations, we set the derivative of this functional to zero giving

$$p(x) = \exp \{ -1 + \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2 \}. \quad (1.108)$$

### Exercise 1.34

The Lagrange multipliers can be found by back substitution of this result into the three constraint equations, leading finally to the result

$$p(x) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\} \quad (1.109)$$

and so the distribution that maximizes the differential entropy is the Gaussian. Note that we did not constrain the distribution to be nonnegative when we maximized the entropy. However, because the resulting distribution is indeed nonnegative, we see with hindsight that such a constraint is not necessary.

### Exercise 1.35

If we evaluate the differential entropy of the Gaussian, we obtain

$$H[x] = \frac{1}{2} \{ 1 + \ln(2\pi\sigma^2) \}. \quad (1.110)$$

Thus we see again that the entropy increases as the distribution becomes broader, i.e., as  $\sigma^2$  increases. This result also shows that the differential entropy, unlike the discrete entropy, can be negative, because  $H(x) < 0$  in (1.110) for  $\sigma^2 < 1/(2\pi e)$ .

Suppose we have a joint distribution  $p(\mathbf{x}, \mathbf{y})$  from which we draw pairs of values of  $\mathbf{x}$  and  $\mathbf{y}$ . If a value of  $\mathbf{x}$  is already known, then the additional information needed to specify the corresponding value of  $\mathbf{y}$  is given by  $-\ln p(\mathbf{y}|\mathbf{x})$ . Thus the average additional information needed to specify  $\mathbf{y}$  can be written as

$$H[\mathbf{y}|\mathbf{x}] = - \iint p(\mathbf{y}, \mathbf{x}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{y} d\mathbf{x} \quad (1.111)$$

**Exercise 1.37**

which is called the *conditional entropy* of  $\mathbf{y}$  given  $\mathbf{x}$ . It is easily seen, using the product rule, that the conditional entropy satisfies the relation

$$H[\mathbf{x}, \mathbf{y}] = H[\mathbf{y}|\mathbf{x}] + H[\mathbf{x}] \quad (1.112)$$

where  $H[\mathbf{x}, \mathbf{y}]$  is the differential entropy of  $p(\mathbf{x}, \mathbf{y})$  and  $H[\mathbf{x}]$  is the differential entropy of the marginal distribution  $p(\mathbf{x})$ . Thus the information needed to describe  $\mathbf{x}$  and  $\mathbf{y}$  is given by the sum of the information needed to describe  $\mathbf{x}$  alone plus the additional information required to specify  $\mathbf{y}$  given  $\mathbf{x}$ .

### 1.6.1 Relative entropy and mutual information

So far in this section, we have introduced a number of concepts from information theory, including the key notion of entropy. We now start to relate these ideas to pattern recognition. Consider some unknown distribution  $p(\mathbf{x})$ , and suppose that we have modelled this using an approximating distribution  $q(\mathbf{x})$ . If we use  $q(\mathbf{x})$  to construct a coding scheme for the purpose of transmitting values of  $\mathbf{x}$  to a receiver, then the average *additional* amount of information (in nats) required to specify the value of  $\mathbf{x}$  (assuming we choose an efficient coding scheme) as a result of using  $q(\mathbf{x})$  instead of the true distribution  $p(\mathbf{x})$  is given by

$$\begin{aligned} \text{KL}(p||q) &= - \int p(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x} - \left( - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \right) \\ &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x}. \end{aligned} \quad (1.113)$$

This is known as the *relative entropy* or *Kullback-Leibler divergence*, or *KL divergence* (Kullback and Leibler, 1951), between the distributions  $p(\mathbf{x})$  and  $q(\mathbf{x})$ . Note that it is not a symmetrical quantity, that is to say  $\text{KL}(p||q) \neq \text{KL}(q||p)$ .

We now show that the Kullback-Leibler divergence satisfies  $\text{KL}(p||q) \geq 0$  with equality if, and only if,  $p(\mathbf{x}) = q(\mathbf{x})$ . To do this we first introduce the concept of *convex* functions. A function  $f(x)$  is said to be convex if it has the property that every chord lies on or above the function, as shown in Figure 1.31. Any value of  $x$  in the interval from  $x = a$  to  $x = b$  can be written in the form  $\lambda a + (1 - \lambda)b$  where  $0 \leq \lambda \leq 1$ . The corresponding point on the chord is given by  $\lambda f(a) + (1 - \lambda)f(b)$ ,



**Claude Shannon**  
1916–2001

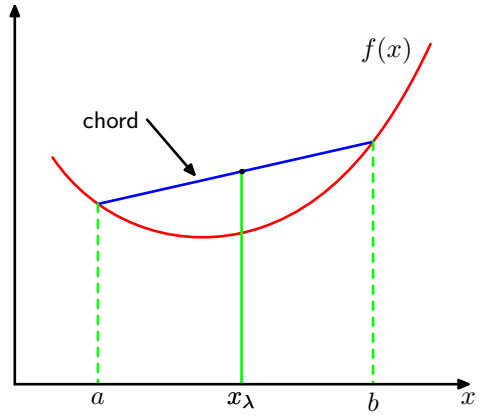
After graduating from Michigan and MIT, Shannon joined the AT&T Bell Telephone laboratories in 1941. His paper 'A Mathematical Theory of Communication' published in the *Bell System Technical Journal* in

1948 laid the foundations for modern information the-

ory. This paper introduced the word 'bit', and his concept that information could be sent as a stream of 1s and 0s paved the way for the communications revolution. It is said that von Neumann recommended to Shannon that he use the term entropy, not only because of its similarity to the quantity used in physics, but also because "nobody knows what entropy really is, so in any discussion you will always have an advantage".



**Figure 1.31** A convex function  $f(x)$  is one for which every chord (shown in blue) lies on or above the function (shown in red).



and the corresponding value of the function is  $f(\lambda a + (1 - \lambda)b)$ . Convexity then implies

$$f(\lambda a + (1 - \lambda)b) \leq \lambda f(a) + (1 - \lambda)f(b). \quad (1.114)$$

This is equivalent to the requirement that the second derivative of the function be everywhere positive. Examples of convex functions are  $x \ln x$  (for  $x > 0$ ) and  $x^2$ . A function is called *strictly convex* if the equality is satisfied only for  $\lambda = 0$  and  $\lambda = 1$ . If a function has the opposite property, namely that every chord lies on or below the function, it is called *concave*, with a corresponding definition for *strictly concave*. If a function  $f(x)$  is convex, then  $-f(x)$  will be concave.

Using the technique of proof by induction, we can show from (1.114) that a convex function  $f(x)$  satisfies

$$f\left(\sum_{i=1}^M \lambda_i x_i\right) \leq \sum_{i=1}^M \lambda_i f(x_i) \quad (1.115)$$

where  $\lambda_i \geq 0$  and  $\sum_i \lambda_i = 1$ , for any set of points  $\{x_i\}$ . The result (1.115) is known as *Jensen's inequality*. If we interpret the  $\lambda_i$  as the probability distribution over a discrete variable  $x$  taking the values  $\{x_i\}$ , then (1.115) can be written

$$f(\mathbb{E}[x]) \leq \mathbb{E}[f(x)] \quad (1.116)$$

where  $\mathbb{E}[\cdot]$  denotes the expectation. For continuous variables, Jensen's inequality takes the form

$$f\left(\int \mathbf{x} p(\mathbf{x}) d\mathbf{x}\right) \leq \int f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}. \quad (1.117)$$

We can apply Jensen's inequality in the form (1.117) to the Kullback-Leibler divergence (1.113) to give

$$\text{KL}(p||q) = - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} \geq - \ln \int q(\mathbf{x}) d\mathbf{x} = 0 \quad (1.118)$$

where we have used the fact that  $-\ln x$  is a convex function, together with the normalization condition  $\int q(\mathbf{x}) d\mathbf{x} = 1$ . In fact,  $-\ln x$  is a strictly convex function, so the equality will hold if, and only if,  $q(\mathbf{x}) = p(\mathbf{x})$  for all  $\mathbf{x}$ . Thus we can interpret the Kullback-Leibler divergence as a measure of the dissimilarity of the two distributions  $p(\mathbf{x})$  and  $q(\mathbf{x})$ .

We see that there is an intimate relationship between data compression and density estimation (i.e., the problem of modelling an unknown probability distribution) because the most efficient compression is achieved when we know the true distribution. If we use a distribution that is different from the true one, then we must necessarily have a less efficient coding, and on average the additional information that must be transmitted is (at least) equal to the Kullback-Leibler divergence between the two distributions.

Suppose that data is being generated from an unknown distribution  $p(\mathbf{x})$  that we wish to model. We can try to approximate this distribution using some parametric distribution  $q(\mathbf{x}|\boldsymbol{\theta})$ , governed by a set of adjustable parameters  $\boldsymbol{\theta}$ , for example a multivariate Gaussian. One way to determine  $\boldsymbol{\theta}$  is to minimize the Kullback-Leibler divergence between  $p(\mathbf{x})$  and  $q(\mathbf{x}|\boldsymbol{\theta})$  with respect to  $\boldsymbol{\theta}$ . We cannot do this directly because we don't know  $p(\mathbf{x})$ . Suppose, however, that we have observed a finite set of training points  $\mathbf{x}_n$ , for  $n = 1, \dots, N$ , drawn from  $p(\mathbf{x})$ . Then the expectation with respect to  $p(\mathbf{x})$  can be approximated by a finite sum over these points, using (1.35), so that

$$\text{KL}(p\|q) \simeq \sum_{n=1}^N \{-\ln q(\mathbf{x}_n|\boldsymbol{\theta}) + \ln p(\mathbf{x}_n)\}. \quad (1.119)$$

The second term on the right-hand side of (1.119) is independent of  $\boldsymbol{\theta}$ , and the first term is the negative log likelihood function for  $\boldsymbol{\theta}$  under the distribution  $q(\mathbf{x}|\boldsymbol{\theta})$  evaluated using the training set. Thus we see that minimizing this Kullback-Leibler divergence is equivalent to maximizing the likelihood function.

Now consider the joint distribution between two sets of variables  $\mathbf{x}$  and  $\mathbf{y}$  given by  $p(\mathbf{x}, \mathbf{y})$ . If the sets of variables are independent, then their joint distribution will factorize into the product of their marginals  $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$ . If the variables are not independent, we can gain some idea of whether they are 'close' to being independent by considering the Kullback-Leibler divergence between the joint distribution and the product of the marginals, given by

$$\begin{aligned} I[\mathbf{x}, \mathbf{y}] &\equiv \text{KL}(p(\mathbf{x}, \mathbf{y})\|p(\mathbf{x})p(\mathbf{y})) \\ &= - \iint p(\mathbf{x}, \mathbf{y}) \ln \left( \frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x} d\mathbf{y} \end{aligned} \quad (1.120)$$

which is called the *mutual information* between the variables  $\mathbf{x}$  and  $\mathbf{y}$ . From the properties of the Kullback-Leibler divergence, we see that  $I(\mathbf{x}, \mathbf{y}) \geq 0$  with equality if, and only if,  $\mathbf{x}$  and  $\mathbf{y}$  are independent. Using the sum and product rules of probability, we see that the mutual information is related to the conditional entropy through

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]. \quad (1.121)$$

Thus we can view the mutual information as the reduction in the uncertainty about  $\mathbf{x}$  by virtue of being told the value of  $\mathbf{y}$  (or vice versa). From a Bayesian perspective, we can view  $p(\mathbf{x})$  as the prior distribution for  $\mathbf{x}$  and  $p(\mathbf{x}|\mathbf{y})$  as the posterior distribution after we have observed new data  $\mathbf{y}$ . The mutual information therefore represents the reduction in uncertainty about  $\mathbf{x}$  as a consequence of the new observation  $\mathbf{y}$ .

## Exercises

- 1.1** (★) **www** Consider the sum-of-squares error function given by (1.2) in which the function  $y(x, \mathbf{w})$  is given by the polynomial (1.1). Show that the coefficients  $\mathbf{w} = \{w_i\}$  that minimize this error function are given by the solution to the following set of linear equations

$$\sum_{j=0}^M A_{ij} w_j = T_i \quad (1.122)$$

where

$$A_{ij} = \sum_{n=1}^N (x_n)^{i+j}, \quad T_i = \sum_{n=1}^N (x_n)^i t_n. \quad (1.123)$$

Here a suffix  $i$  or  $j$  denotes the index of a component, whereas  $(x)^i$  denotes  $x$  raised to the power of  $i$ .

- 1.2** (★) Write down the set of coupled linear equations, analogous to (1.122), satisfied by the coefficients  $w_i$  which minimize the regularized sum-of-squares error function given by (1.4).
- 1.3** (★★) Suppose that we have three coloured boxes  $r$  (red),  $b$  (blue), and  $g$  (green). Box  $r$  contains 3 apples, 4 oranges, and 3 limes, box  $b$  contains 1 apple, 1 orange, and 0 limes, and box  $g$  contains 3 apples, 3 oranges, and 4 limes. If a box is chosen at random with probabilities  $p(r) = 0.2$ ,  $p(b) = 0.2$ ,  $p(g) = 0.6$ , and a piece of fruit is removed from the box (with equal probability of selecting any of the items in the box), then what is the probability of selecting an apple? If we observe that the selected fruit is in fact an orange, what is the probability that it came from the green box?
- 1.4** (★★) **www** Consider a probability density  $p_x(x)$  defined over a continuous variable  $x$ , and suppose that we make a nonlinear change of variable using  $x = g(y)$ , so that the density transforms according to (1.27). By differentiating (1.27), show that the location  $\hat{y}$  of the maximum of the density in  $y$  is not in general related to the location  $\hat{x}$  of the maximum of the density over  $x$  by the simple functional relation  $\hat{x} = g(\hat{y})$  as a consequence of the Jacobian factor. This shows that the maximum of a probability density (in contrast to a simple function) is dependent on the choice of variable. Verify that, in the case of a linear transformation, the location of the maximum transforms in the same way as the variable itself.

- 1.5** (★) Using the definition (1.38) show that  $\text{var}[f(x)]$  satisfies (1.39).