

where we have used $-\ln x$ 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 $\boldsymbol{\theta}$. 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 do not 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 (2.40), so that

$$\text{KL}(p||q) \simeq \frac{1}{N} \sum_{n=1}^N \left\{ -\ln q(\mathbf{x}_n|\boldsymbol{\theta}) + \ln p(\mathbf{x}_n) \right\}. \quad (2.106)$$

The second term on the right-hand side of (2.106) 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 log likelihood function.

Exercise 2.34

2.5.6 Conditional entropy

Now consider the joint distribution between two sets of variables \mathbf{x} and \mathbf{y} given by $p(\mathbf{x}, \mathbf{y})$ 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

$$H[\mathbf{y}|\mathbf{x}] = - \iint p(\mathbf{y}, \mathbf{x}) \ln p(\mathbf{y}|\mathbf{x}) d\mathbf{y} d\mathbf{x}, \quad (2.107)$$

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:

Exercise 2.35

$$H[\mathbf{x}, \mathbf{y}] = H[\mathbf{y}|\mathbf{x}] + H[\mathbf{x}] \quad (2.108)$$

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} .