# 10. Information theory

# 10.1. Shannon information and entropy

10.2. Kullback Leibler divergence

10.3. Flow of information

# 10.1. Shannon information and entropy

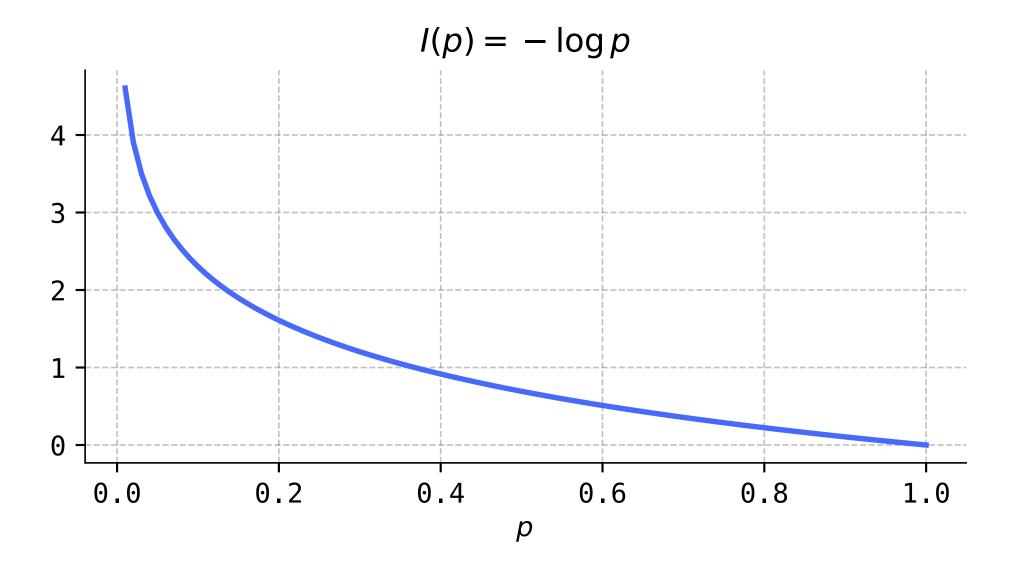
#### Definition 10.1

Let P be a probability measure on a finite sample space S with mass function p(s). The (Shannon) information content of the sample point  $s \in S$ , denoted  $I_P(s)$ , is defined to be

$$I_P(s) \stackrel{ ext{def}}{=} -\log(p(s)).$$

The information content is also called the *surprisal*.

If the probability measure P is clear from context, we will write I(s) in place of  $I_P(s)$ . If  $\mathbf{X}$  is a random vector with finite range and probability measure  $P_{\mathbf{X}}$ , we will write  $I_{\mathbf{X}}(\mathbf{x})$  in place of  $I_{P_{\mathbf{X}}}(\mathbf{x})$ .



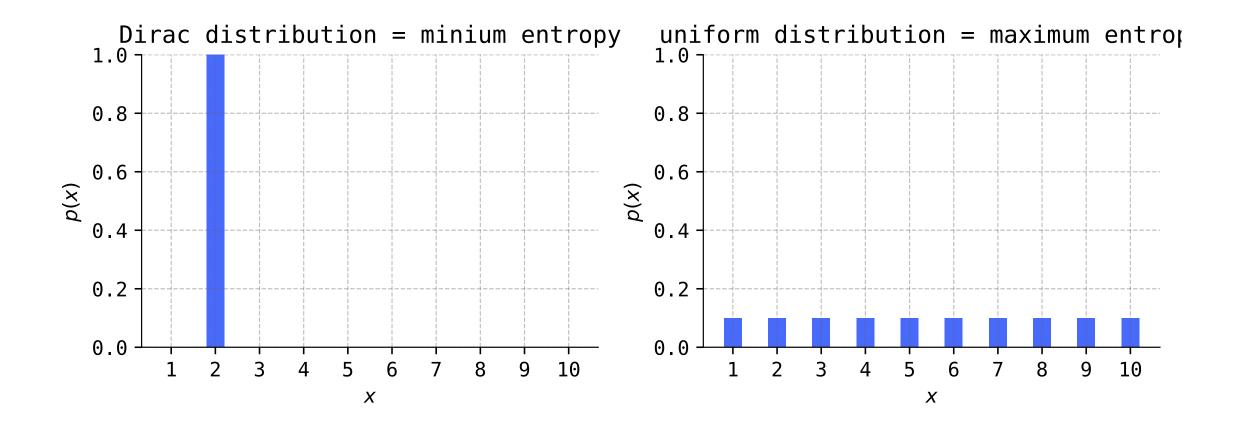
#### Definition 10.2

Let P be a probability measure on a finite sample space S with mass function p(s). The (Shannon) entropy of P, denoted H(P), is defined to be

$$H(P) \stackrel{\mathrm{def}}{=} \sum_{s \in S} p(s) I_P(s).$$

The entropy is also called the *uncertainty*.

If  ${\bf X}$  is a random vector with finite range and probability measure  $P_{\bf X}$ , we will write  $H({\bf X})$  in place of  $H(P_{\bf X})$ . If we write the vector in terms of its component random variables  ${\bf X}=(X_1,\ldots,X_m)$ , then we shall also write  $H(X_1,\ldots,X_m)$  in place of  $H(P_{\bf X})$  and call this the *joint entropy* of the random variables  $X_1,\ldots,X_m$ .





Do problems 1 and 2 on the worksheet.

#### **Definition 10.3**

Let P and Q be two probability measures on a finite sample space S with mass functions p(s) and q(s). Suppose they satisfy the following condition:

• Absolute continuity. For all  $s \in S$ , if q(s) = 0, then p(s) = 0. Or equivalently, the support of q(s) contains the support of p(s).

Then the *cross entropy* from P to Q, denoted  $H_P(Q)$ , is defined by

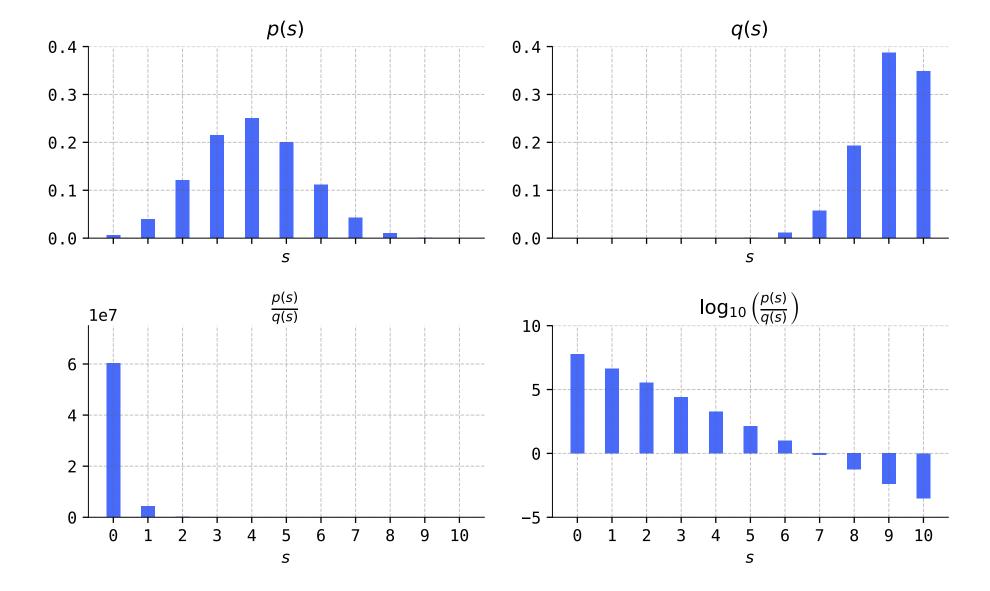
$$H_P(Q) \stackrel{ ext{def}}{=} E_{s \sim p(s)} \left[ I_Q(s) 
ight] = - \sum_{s \in S} p(s) \log(q(s)).$$

As usual, if  $P_{\mathbf{X}}$  and  $P_{\mathbf{Y}}$  are the probability measures of two random vectors  $\mathbf{X}$  and  $\mathbf{Y}$  with finite ranges, we will write  $H_{\mathbf{Y}}(\mathbf{X})$  in place of  $H_{P_{\mathbf{Y}}}(P_{\mathbf{X}})$ .



Do problem 3 on the worksheet.

# 10.2. Kullback Leibler divergence



Let P and Q be two probability measures on a finite sample space S with mass functions p(s) and q(s). Suppose they satisfy the following condition:

• Absolute continuity. For all  $s \in S$ , if q(s) = 0, then p(s) = 0. Or equivalently, the support of q(s) contains the support of p(s).

Then the *Kullback-Leibler divergence* (or just *KL divergence*) from P to Q, denoted  $D(P \parallel Q)$ , is the mean order of relative magnitude of P to Q. Precisely, it is given by

$$D(P \parallel Q) \stackrel{ ext{def}}{=} E_{s \sim p(s)} \left[ \log \left( rac{p(s)}{q(s)} 
ight) 
ight] = \sum_{s \in S} p(s) \log \left( rac{p(s)}{q(s)} 
ight).$$

The KL divergence is also called the *relative entropy*.

As always, if  $P_{\mathbf{X}}$  and  $P_{\mathbf{Y}}$  are the probability measures of two random vectors  $\mathbf{X}$  and  $\mathbf{Y}$  with finite ranges, we will write  $D(\mathbf{Y} \parallel \mathbf{X})$  in place of  $D(P_{\mathbf{Y}} \parallel P_{\mathbf{X}})$ .



To problem 4 on the worksheet.

#### Theorem 10.1 (KL divergence and entropy)

Let P and Q be two probability measures on a finite sample space S. Then

$$D(P \parallel Q) = H_P(Q) - H(P).$$

#### Theorem 10.3 (Gibbs' inequality)

Let P and Q be two probability measures on a finite probability space S satisfying the absolute continuity condition in <u>Definition 10.4</u>. Then

$$D(P \parallel Q) \geq 0$$
,

with equality if and only if P=Q.

## Corollary 10.1 (Uniform distributions maximize entropy)

Let P be a probability measures on a finite sample space S. Then

$$H(P) \leq \log |S|,$$

with equality if and only if P is uniform.



Do problem 5 on the worksheet.