

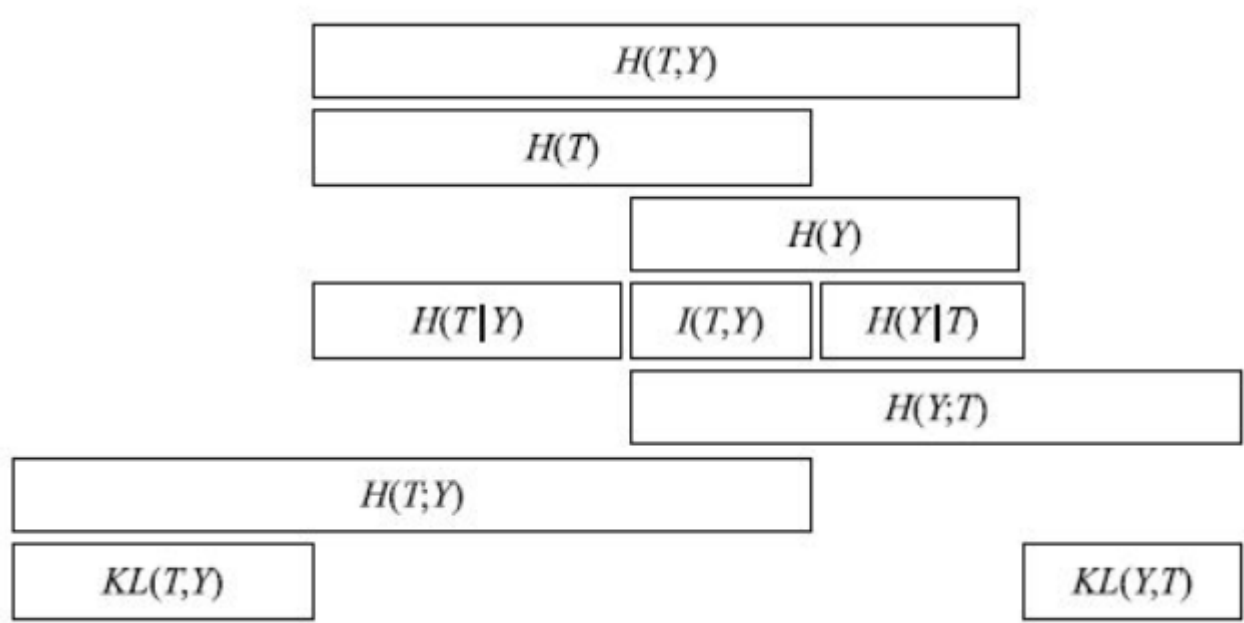
信息论基本概念汇总

总览:

Name	Formula	(Dis)similarity	(A)symmetry
Joint Information	$H(T, Y) = - \sum_t \sum_y p(t, y) \log_2 p(t, y)$	Inapplicable	Symmetry
Mutual Information	$I(T, Y) = \sum_t \sum_y p(t, y) \log_2 \frac{p(t, y)}{p(t)p(y)}$	Similarity	Symmetry
Conditional Entropy	$H(Y T) = - \sum_t \sum_y p(t, y) \log_2 p(y t)$	Dissimilarity	Asymmetry
Cross Entropy	$H(T; Y) = - \sum_z p_t(z) \log_2 p_y(z)$	Dissimilarity	Asymmetry
KL Divergence	$KL(T, Y) = \sum_z p_t(z) \log_2 \frac{p_t(z)}{p_y(z)}$	Dissimilarity	Asymmetry

这些度量中，互信息可以用来衡量相似性（越大越相似），而条件熵、交叉熵和相对熵（即KL散度）可以用来度量相异性（越大越不相似）。

这些量之间的关系如下



1. 信息量

如果事件 X 发生，那么 $p(x)$ 能为“事件 x 发生”所提供的信息量：

$$h(X) = -\log_2 p(x)$$

也就是消除事情不确定性所需要的信息量，单位是比特。

2. 熵

熵是接收的每条消息中包含的信息的平均量，它是不确定性的度量，越随机的信号源其熵越大。

离散：

$$\mathbb{H}(X) = - \sum_x p(x_i) \log_2 p(x_i)$$

连续：

$$\mathbb{H}(X) = - \int p(x) \log_2 p(x)$$

在最优化理论中，很多算法用熵作为优化目标，Watanabe也提出过“学习就是一个熵减的过程”，算法学习的过程就是信息不确定性减小的过程。

3. 联合熵

度量二维随机变量的不确定性：

$$\mathbb{H}(X, Y) = - \sum_i \sum_j p(x_i, y_j) \log_2 p(x_i, y_j)$$

4. 条件熵

$\mathbb{H}(Y|X)$ 表示已知 X ，求 Y 的平均不确定性

$$\mathbb{H}(Y|X) = - \sum_i \sum_j p(x_i, y_j) \log_2 p(y_j|x_i)$$

$$\mathbb{H}(Y|X) = \sum_i p(x_i) \mathbb{H}(Y|x_i)$$

推导过程如下：

$$\begin{aligned} \mathbb{H}(Y|X) &= - \sum_i \sum_j p(x_i) p(y_i|x_i) \log_2 p(y_i|x_i) \\ &= - \sum_i p(x_i) \sum_j p(y_i|x_i) \log_2 p(y_i|x_i) \\ &= \sum_i p(x_i) \mathbb{H}(Y|x_i) \end{aligned}$$

由联合熵和条件熵可得：

$$\begin{aligned}
\mathbb{H}(X, Y) &= - \sum_i \sum_j p(x_i, y_j) \log_2 p(x_i, y_j) \\
&= - \sum_i \sum_j p(x_i, y_j) \log_2 p(y_j | x_i) \\
&\quad + \sum_i \left(\sum_j p(x_i, y_j) \right) \log_2 p(x_i) \\
&= \mathbb{H}(Y|X) + \mathbb{H}(X)
\end{aligned}$$

5. 相对熵

又称为KL散度。主要用来衡量两个分布的相似度。假设连续随机变量 x ，真实的概率分布为 $p(x)$ ，模型得到的近似分布为 $q(x)$ 。

离散：

$$\begin{aligned}
\mathbb{KL}(p||q) &= - \sum_i p(x_i) \ln q(x_i) - \left(- \sum_i p(x_i) \ln p(x_i) \right) \\
&= \sum_i p(x_i) \ln \frac{p(x_i)}{q(x_i)}
\end{aligned}$$

连续：

$$\begin{aligned}
\mathbb{KL}(p||q) &= - \int_x p(x) \ln p(x) + p(x) \ln q(x) \\
&= \int_x p(x) \ln \frac{p(x)}{q(x)}
\end{aligned}$$

对离散变量的相对熵：

$$\begin{aligned}
\mathbb{KL}(p||q) &= - \sum_i p(x_i) \ln q(x_i) - \left(- \sum_i p(x_i) \ln p(x_i) \right) \\
&= \mathbb{H}(p, q) - \mathbb{H}(p)
\end{aligned}$$

其中 $\mathbb{H}(p, q)$ 被称为交叉熵。

6. 交叉熵

用来衡量两个分布之间的相似程度， p 为真实概率分布， q 为模型预测的概率分布：

$$\mathbb{H}(p, q) = - \sum_i p(x_i) \ln q(x_i)$$

7. 互信息

相对熵是衡量同一个变量的两个一维分布之间的相似性，而互信息是用来衡量两个相同的一维分布变量之间的独立性。

互信息 $\mathbb{I}(p, q)$ 是衡量联合分布 $p(x, y)$ 和 $p(x)p(y)$ 分布之间的关系，即它们之间的相关系数

$$\begin{aligned}
\mathbb{I}(X, Y) &= \mathbb{KL}(p(x, y) \| p(x)p(y)) \\
&= \sum_i \sum_j p(x_i, y_j) \ln \frac{p(x_i, y_j)}{p(x_i)p(y_j)} \\
&= -\mathbb{H}(X, Y) + \mathbb{H}(X) + \mathbb{H}(Y) \\
&= \mathbb{H}(X) - \mathbb{H}(X|Y) \\
&= \mathbb{H}(Y) - \mathbb{H}(Y|X)
\end{aligned}$$

$\mathbb{I}(X, Y)$ 反映的是在知道了 Y 的值以后 X 的不确定性的减少量。 可以理解为 Y 的值透露了多少关于 X 的信息量。

实际上，互信息体现了两变量之间的依赖程度：如果 $\mathbb{I}(X, Y) \gg 0$ ，表明 X 和 Y 是高度相关的；如果 $\mathbb{I}(X, Y) = 0$ ，表明 X 和 Y 是相互独立的；如果 $\mathbb{I}(X, Y) \ll 0$ ，表明 Y 的出现不但未使 X 得不确定性减少，反而还增大了 X 的不确定性，常常是不利的。