

# Cross entropy

In [information theory](#), the **cross entropy** between two [probability distributions](#) ***p*** and ***q*** over the same underlying set of events measures the average number of [bits](#) needed to identify an event drawn from the set if a coding scheme used for the set is optimized for an estimated probability distribution ***q***, rather than the true distribution ***p***.

## Contents

- Definition
- Motivation
- Estimation
- Relation to log-likelihood
- Cross-entropy minimization
- Cross-entropy loss function and logistic regression
- See also
- References
- External links

## Definition

The cross entropy of the distribution ***q*** relative to a distribution ***p*** over a given set is defined as follows:

$$H(p, q) = - \mathbf{E}_p[\log q].$$

The definition may be formulated using the [Kullback–Leibler divergence](#)  $D_{\text{KL}}(p\|q)$  from ***p*** of ***q*** (also known as the *relative entropy* of ***q*** with respect to ***p***).

$$H(p, q) = H(p) + D_{\text{KL}}(p\|q),$$

where  $H(p)$  is the [entropy](#) of ***p***.

For [discrete probability distributions](#) ***p*** and ***q*** with the same [support](#)  $\mathcal{X}$  this means

$$H(p, q) = - \sum_{x \in \mathcal{X}} p(x) \log q(x)$$

(Eq.1)

The situation for [continuous distributions](#) is analogous. We have to assume that ***p*** and ***q*** are [absolutely continuous](#) with respect to some reference measure ***r*** (usually ***r*** is a [Lebesgue measure](#) on a [Borel σ-algebra](#)). Let ***P*** and ***Q*** be probability density functions of ***p*** and ***q*** with respect to ***r***. Then

$$- \int_{\mathcal{X}} P(x) \log Q(x) \, d\mathbf{r}(x) = \mathbf{E}_p[-\log Q]$$

and therefore

$$H(p, q) = - \int_{\mathcal{X}} P(x) \log Q(x) \, d\mathbf{r}(x)$$

(Eq.2)

NB: The notation  $H(p, q)$  is also used for a different concept, the [joint entropy](#) of ***p*** and ***q***.

## Motivation

In [information theory](#), the [Kraft–McMillan theorem](#) establishes that any directly decodable coding scheme for coding a message to identify one value ***x<sub>i</sub>*** out of a set of possibilities  $\{x_1, \dots, x_n\}$  can be seen as representing an implicit probability distribution  $q(x_i) = \left(\frac{1}{2}\right)^{l_i}$  over  $\{x_1, \dots, x_n\}$ , where  $l_i$  is the length of the code for ***x<sub>i</sub>*** in bits. Therefore, cross entropy can be interpreted as the expected message-length per

datum when a wrong distribution  $\mathbf{q}$  is assumed while the data actually follows a distribution  $\mathbf{p}$ . That is why the expectation is taken over the true probability distribution  $\mathbf{p}$  and not  $\mathbf{q}$ . Indeed the expected message-length under the true distribution  $\mathbf{p}$  is,

$$\mathbb{E}_{\mathbf{p}}[l] = -\mathbb{E}_{\mathbf{p}}\left[\frac{\ln q(\mathbf{x})}{\ln(2)}\right] = -\mathbb{E}_{\mathbf{p}}[\log_2 q(\mathbf{x})] = -\sum_{\mathbf{x}_i} \mathbf{p}(\mathbf{x}_i) \log_2 q(\mathbf{x}_i) = -\sum_{\mathbf{x}} \mathbf{p}(\mathbf{x}) \log_2 q(\mathbf{x}) = H(\mathbf{p}, \mathbf{q})$$

## Estimation

There are many situations where cross-entropy needs to be measured but the distribution of  $\mathbf{p}$  is unknown. An example is [language modeling](#), where a model is created based on a training set  $\mathbf{T}$ , and then its cross-entropy is measured on a test set to assess how accurate the model is in predicting the test data. In this example,  $\mathbf{p}$  is the true distribution of words in any corpus, and  $\mathbf{q}$  is the distribution of words as predicted by the model. Since the true distribution is unknown, cross-entropy cannot be directly calculated. In these cases, an estimate of cross-entropy is calculated using the following formula:

$$H(\mathbf{T}, \mathbf{q}) = -\sum_{i=1}^N \frac{1}{N} \log_2 q(\mathbf{x}_i)$$

where  $N$  is the size of the test set, and  $q(\mathbf{x})$  is the probability of event  $\mathbf{x}$  estimated from the training set. The sum is calculated over  $N$ . This is a Monte Carlo estimate of the true cross entropy, where the test set is treated as samples from  $\mathbf{p}(\mathbf{x})$ .

## Relation to log-likelihood

In classification problems we want to estimate the probability of different outcomes. If the estimated probability of outcome  $i$  is  $q_i$ , while the frequency (empirical probability) of outcome  $i$  in the training set is  $p_i$ , and there are  $N$  samples in the training set, then the likelihood of the training set is proportional to

$$\prod_i q_i^{Np_i}$$

so the log-likelihood, divided by  $N$  is

$$\frac{1}{N} \log \prod_i q_i^{Np_i} = \sum_i p_i \log q_i = -H(\mathbf{p}, \mathbf{q})$$

so that maximizing the likelihood is the same as minimizing the cross entropy.

## Cross-entropy minimization

Cross-entropy minimization is frequently used in optimization and rare-event probability estimation; see the [cross-entropy method](#).

When comparing a distribution  $\mathbf{q}$  against a fixed reference distribution  $\mathbf{p}$ , cross entropy and KL divergence are identical up to an additive constant (since  $\mathbf{p}$  is fixed): both take on their minimal values when  $\mathbf{p} = \mathbf{q}$ , which is  $\mathbf{0}$  for KL divergence, and  $\mathbf{H}(\mathbf{p})$  for cross entropy.<sup>[1]</sup> In the engineering literature, the principle of minimising KL Divergence (Kullback's "[Principle of Minimum Discrimination Information](#)") is often called the **Principle of Minimum Cross-Entropy** (MCE), or **Minxent**.

However, as discussed in the article *Kullback–Leibler divergence*, sometimes the distribution  $\mathbf{q}$  is the fixed prior reference distribution, and the distribution  $\mathbf{p}$  is optimised to be as close to  $\mathbf{q}$  as possible, subject to some constraint. In this case the two minimisations are *not* equivalent. This has led to some ambiguity in the literature, with some authors attempting to resolve the inconsistency by redefining cross-entropy to be  $D_{\text{KL}}(\mathbf{p}||\mathbf{q})$ , rather than  $H(\mathbf{p}, \mathbf{q})$ .

## Cross-entropy loss function and logistic regression

Cross entropy can be used to define a loss function in [machine learning](#) and [optimization](#). The true probability  $p_i$  is the true label, and the given distribution  $q_i$  is the predicted value of the current model.

More specifically, consider [logistic regression](#), which (among other things) can be used to classify observations into two possible classes (often simply labelled  $\mathbf{0}$  and  $\mathbf{1}$ ). The output of the model for a given observation, given a vector of input features  $\mathbf{x}$ , can be interpreted as a probability, which serves as the basis for classifying the observation. The probability is modeled using the [logistic function](#)  $g(\mathbf{z}) = 1/(1 + e^{-\mathbf{z}})$  where  $\mathbf{z}$  is some function of the input vector  $\mathbf{x}$ , commonly just a linear function. The probability of the output  $\mathbf{y} = \mathbf{1}$  is given by

$$q_{y=1} = \hat{y} \equiv g(\mathbf{w} \cdot \mathbf{x}) = 1/(1 + e^{-\mathbf{w} \cdot \mathbf{x}}),$$

where the vector of weights  $\mathbf{w}$  is optimized through some appropriate algorithm such as gradient descent. Similarly, the complementary probability of finding the output  $y = 0$  is simply given by

$$q_{y=0} = 1 - \hat{y}$$

Having set up our notation,  $p \in \{y, 1 - y\}$  and  $q \in \{\hat{y}, 1 - \hat{y}\}$ , we can use cross entropy to get a measure of dissimilarity between  $p$  and  $q$ :

$$H(p, q) = - \sum_i p_i \log q_i = - y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

Logistic regression typically optimizes the log loss for all the observations on which it is trained, which is the same as optimizing the average cross-entropy in the sample. For example, suppose we have  $N$  samples with each sample indexed by  $n = 1, \dots, N$ . The average of the loss function is then given by:

$$J(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N H(p_n, q_n) = - \frac{1}{N} \sum_{n=1}^N \left[ y_n \log \hat{y}_n + (1 - y_n) \log(1 - \hat{y}_n) \right],$$

where  $\hat{y}_n \equiv g(\mathbf{w} \cdot \mathbf{x}_n) = 1/(1 + e^{-\mathbf{w} \cdot \mathbf{x}_n})$ , with  $g(z)$  the logistic function as before.

The logistic loss is sometimes called cross-entropy loss. It is also known as log loss (In this case, the binary label is often denoted by  $\{-1, +1\}$ ).<sup>[2]</sup>

## See also

- Cross-entropy method
- Logistic regression
- Conditional entropy
- Maximum likelihood estimation
- Mutual information

## References

- Ian Goodfellow, Yoshua Bengio, and Aaron Courville (2016). Deep Learning. MIT Press. Online (<http://www.deeplearningbook.org>)
- Murphy, Kevin (2012). *Machine Learning: A Probabilistic Perspective*. MIT. ISBN 978-0262018029.

## External links

- What is cross-entropy, and why use it? (<http://www.cse.unsw.edu.au/~billw/cs9444/crossentropy.html>)
- Cross Entropy (<http://heliosphan.org/cross-entropy.html>)

Retrieved from "[https://en.wikipedia.org/w/index.php?title=Cross\\_entropy&oldid=941508274](https://en.wikipedia.org/w/index.php?title=Cross_entropy&oldid=941508274)"

This page was last edited on 19 February 2020, at 00:23 (UTC).

Text is available under the Creative Commons Attribution-ShareAlike License; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy. Wikipedia® is a registered trademark of the Wikimedia Foundation, Inc., a non-profit organization.