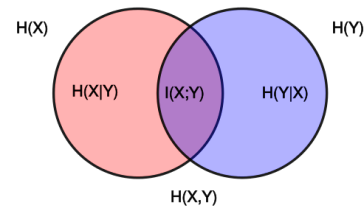


Conditional entropy

In information theory, the **conditional entropy** (or **equivocation**) quantifies the amount of information needed to describe the outcome of a random variable **Y** given that the value of another random variable **X** is known. Here, information is measured in shannons, nats, or hartleys. The *entropy of Y conditioned on X* is written as **H(Y|X)**.



Venn diagram showing additive and subtractive relationships various information measures associated with correlated variables **X** and **Y**. The area contained by both circles is the joint entropy **H(X, Y)**. The circle on the left (red and violet) is the individual entropy **H(X)**, with the red being the conditional entropy **H(X|Y)**. The circle on the right (blue and violet) is **H(Y)**, with the blue being **H(Y|X)**. The violet is the mutual information **I(X; Y)**.

Contents

- Definition
- Motivation
- Properties
 - Conditional entropy equals zero
 - Conditional entropy of independent random variables
 - Chain rule
 - Bayes' rule
 - Other properties
- Conditional differential entropy
 - Definition
 - Properties
 - Relation to estimator error
- Generalization to quantum theory
- See also
- References

Definition

The conditional entropy of **Y** given **X** is defined as

$$H(Y|X) = - \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x)}$$

(Eq.1)

where **X** and **Y** denote the support sets of **X** and **Y**.

Note: It is conventioned that the expressions **0 log 0** and **0 log c/0** for fixed **c > 0** should be treated as being equal to zero. This is because $\lim_{\theta \rightarrow 0^+} \theta \log c/\theta = 0$ and $\lim_{\theta \rightarrow 0^+} \theta \log \theta = 0^{[1]}$

Motivation

Let **H(Y|X = x)** be the entropy of the discrete random variable **Y** conditioned on the discrete random variable **X** taking a certain value **x**. Denote the support sets of **X** and **Y** by **X** and **Y**. Let **Y** have probability mass function **p_Y(y)**. The unconditional entropy of **Y** is calculated as **H(Y) := E[I(Y)]**, i.e.

$$H(Y) = \sum_{y \in \mathcal{Y}} \Pr(Y = y) I(y) = - \sum_{y \in \mathcal{Y}} p_Y(y) \log_2 p_Y(y),$$

where **I(y_i)** is the information content of the outcome of **Y** taking the value **y_i**. The entropy of **Y** conditioned on **X** taking the value **x** is defined analogously by conditional expectation:

$$H(Y|X = x) = - \sum_{y \in \mathcal{Y}} \Pr(Y = y|X = x) \log_2 \Pr(Y = y|X = x).$$

H(Y|X) is the result of averaging **H(Y|X = x)** over all possible values **x** that **X** may take.

Given discrete random variables **X** with image **X** and **Y** with image **Y**, the conditional entropy of **Y** given **X** is defined as the weighted sum of **H(Y|X = x)** for each possible value of **x**, using **p(x)** as the weights:^{[2]:15}

$$\begin{aligned}
H(Y|X) &\equiv \sum_{x \in \mathcal{X}} p(x) H(Y|X=x) \\
&= - \sum_{x \in \mathcal{X}} p(x) \sum_{y \in \mathcal{Y}} p(y|x) \log p(y|x) \\
&= - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(y|x) \\
&= - \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log p(y|x) \\
&= - \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x)} \\
&= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log \frac{p(x)}{p(x, y)}.
\end{aligned}$$

Properties

Conditional entropy equals zero

$H(Y|X) = 0$ if and only if the value of Y is completely determined by the value of X .

Conditional entropy of independent random variables

Conversely, $H(Y|X) = H(Y)$ if and only if Y and X are independent random variables.

Chain rule

Assume that the combined system determined by two random variables X and Y has joint entropy $H(X, Y)$, that is, we need $H(X, Y)$ bits of information on average to describe its exact state. Now if we first learn the value of X , we have gained $H(X)$ bits of information. Once X is known, we only need $H(X, Y) - H(X)$ bits to describe the state of the whole system. This quantity is exactly $H(Y|X)$, which gives the *chain rule* of conditional entropy:

$$H(Y|X) = H(X, Y) - H(X).^{[2]:17}$$

The chain rule follows from the above definition of conditional entropy:

$$\begin{aligned}
H(Y|X) &= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log \left(\frac{p(x)}{p(x, y)} \right) \\
&= - \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log(p(x, y)) + \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log(p(x)) \\
&= H(X, Y) + \sum_{x \in \mathcal{X}} p(x) \log(p(x)) \\
&= H(X, Y) - H(X).
\end{aligned}$$

In general, a chain rule for multiple random variables holds:

$$H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i | X_1, \dots, X_{i-1})^{[2]:22}$$

It has a similar form to chain rule in probability theory, except that addition instead of multiplication is used.

Bayes' rule

Bayes' rule for conditional entropy states

$$H(Y|X) = H(X|Y) - H(X) + H(Y).$$

Proof. $H(Y|X) = H(X, Y) - H(X)$ and $H(X|Y) = H(Y, X) - H(Y)$. Symmetry entails $H(X, Y) = H(Y, X)$. Subtracting the two equations implies Bayes' rule.

If Y is conditionally independent of Z given X we have:

$$H(Y|X, Z) = H(Y|X).$$

Other properties

For any \mathbf{X} and \mathbf{Y} :

$$\begin{aligned}\mathbf{H}(\mathbf{Y}|\mathbf{X}) &\leq \mathbf{H}(\mathbf{Y}) \\ \mathbf{H}(\mathbf{X}, \mathbf{Y}) &= \mathbf{H}(\mathbf{X}|\mathbf{Y}) + \mathbf{H}(\mathbf{Y}|\mathbf{X}) + \mathbf{I}(\mathbf{X}; \mathbf{Y}), \\ \mathbf{H}(\mathbf{X}, \mathbf{Y}) &= \mathbf{H}(\mathbf{X}) + \mathbf{H}(\mathbf{Y}) - \mathbf{I}(\mathbf{X}; \mathbf{Y}), \\ \mathbf{I}(\mathbf{X}; \mathbf{Y}) &\leq \mathbf{H}(\mathbf{X}),\end{aligned}$$

where $\mathbf{I}(\mathbf{X}; \mathbf{Y})$ is the mutual information between \mathbf{X} and \mathbf{Y} .

For independent \mathbf{X} and \mathbf{Y} :

$$\mathbf{H}(\mathbf{Y}|\mathbf{X}) = \mathbf{H}(\mathbf{Y}) \text{ and } \mathbf{H}(\mathbf{X}|\mathbf{Y}) = \mathbf{H}(\mathbf{X})$$

Although the specific-conditional entropy $\mathbf{H}(\mathbf{X}|\mathbf{Y} = \mathbf{y})$ can be either less or greater than $\mathbf{H}(\mathbf{X})$ for a given random variate \mathbf{y} of \mathbf{Y} , $\mathbf{H}(\mathbf{X}|\mathbf{Y})$ can never exceed $\mathbf{H}(\mathbf{X})$.

Conditional differential entropy

Definition

The above definition is for discrete random variables. The continuous version of discrete conditional entropy is called *conditional differential (or continuous) entropy*. Let \mathbf{X} and \mathbf{Y} be a continuous random variables with a joint probability density function $f(\mathbf{x}, \mathbf{y})$. The differential conditional entropy $h(\mathbf{X}|\mathbf{Y})$ is defined as^{[2]:249}

$$h(\mathbf{X}|\mathbf{Y}) = - \int_{\mathbf{x}, \mathbf{y}} f(\mathbf{x}, \mathbf{y}) \log f(\mathbf{x}|\mathbf{y}) \, d\mathbf{x} d\mathbf{y} \quad (\text{Eq.2})$$

Properties

In contrast to the conditional entropy for discrete random variables, the conditional differential entropy may be negative.

As in the discrete case there is a chain rule for differential entropy:

$$h(\mathbf{Y}|\mathbf{X}) = h(\mathbf{X}, \mathbf{Y}) - h(\mathbf{X})$$
^{[2]:253}

Notice however that this rule may not be true if the involved differential entropies do not exist or are infinite.

Joint differential entropy is also used in the definition of the mutual information between continuous random variables:

$$\mathbf{I}(\mathbf{X}, \mathbf{Y}) = h(\mathbf{X}) - h(\mathbf{X}|\mathbf{Y}) = h(\mathbf{Y}) - h(\mathbf{Y}|\mathbf{X})$$

$h(\mathbf{X}|\mathbf{Y}) \leq h(\mathbf{X})$ with equality if and only if \mathbf{X} and \mathbf{Y} are independent.^{[2]:253}

Relation to estimator error

The conditional differential entropy yields a lower bound on the expected squared error of an estimator. For any random variable \mathbf{X} , observation \mathbf{Y} and estimator $\widehat{\mathbf{X}}$ the following holds:^{[2]:255}

$$\mathbb{E} \left[(\mathbf{X} - \widehat{\mathbf{X}}(\mathbf{Y}))^2 \right] \geq \frac{1}{2\pi e} e^{2h(\mathbf{X}|\mathbf{Y})}$$

This is related to the uncertainty principle from quantum mechanics.

Generalization to quantum theory

In quantum information theory, the conditional entropy is generalized to the conditional quantum entropy. The latter can take negative values, unlike its classical counterpart.

See also

- Entropy (information theory)
- Mutual information
- Conditional quantum entropy

- [Variation of information](#)
- [Entropy power inequality](#)
- [Likelihood function](#)

References

1. "David MacKay: Information Theory, Pattern Recognition and Neural Networks: The Book" (<http://www.inference.org.uk/mackay/itprnn/book.html>). *www.inference.org.uk*. Retrieved 2019-10-25.
 2. T. Cover; J. Thomas (1991). *Elements of Information Theory* (<https://archive.org/details/elementsofinform0000cove>). ISBN 0-471-06259-6.
-

Retrieved from "https://en.wikipedia.org/w/index.php?title=Conditional_entropy&oldid=933858139"

This page was last edited on 3 January 2020, at 11:33 (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.