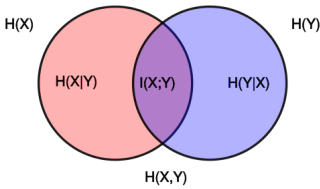


Mutual information

In probability theory and information theory, the **mutual information (MI)** of two random variables is a measure of the mutual dependence between the two variables. More specifically, it quantifies the "amount of information" (in units such as shannons, commonly called bits) obtained about one random variable through observing the other random variable. The concept of mutual information is intricately linked to that of entropy of a random variable, a fundamental notion in information theory that quantifies the expected "amount of information" held in a random variable.

Not limited to real-valued random variables and linear dependence like the correlation coefficient, MI is more general and determines how different the joint distribution of the pair **(*X*, *Y*)** is to the product of the marginal distributions of ***X*** and ***Y***. MI is the expected value of the pointwise mutual information (PMI).

Mutual Information is also known as information gain.



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

In terms of PMFs for discrete distributions

In terms of PDFs for continuous distributions

Motivation

Relation to other quantities

- Nonnegativity
- Symmetry
- Relation to conditional and joint entropy
- Relation to Kullback–Leibler divergence
- Bayesian estimation of mutual information
- Independence assumptions

Variations

- Metric
- Conditional mutual information
- Multivariate mutual information
 - Multivariate statistical independence
- Applications
- Directed information
- Normalized variants
- Weighted variants
- Adjusted mutual information
- Absolute mutual information
- Linear correlation
- For discrete data

Applications

See also

Notes

References

Definition

Let **(*X*, *Y*)** be a pair of random variables with values over the space ***X* × *Y***. If their joint distribution is ***P*(*x*, *y*)** and the marginal distributions are ***P*_{*X*}** and ***P*_{*Y*}**, the mutual information is defined as

$$I(X;Y) = D_{\text{KL}}(P_{(X,Y)} \| P_X \otimes P_Y)$$

where ***D*_{KL}** is the Kullback–Leibler divergence. Notice, as per property of the Kullback–Leibler divergence, that ***I*(*X*; *Y*)** is equal to zero precisely when the joint distribution coincides with the product of the marginals, i.e. when ***X*** and ***Y*** are independent (and hence observing ***Y*** tells your nothing about ***X***). In general ***I*(*X*; *Y*)** is non-negative, it is a measure of the price for encoding **(*X*, *Y*)** as a pair of independent random variables, when in reality they are not.

In terms of PMFs for discrete distributions

The mutual information of two jointly discrete random variables ***X*** and ***Y*** is calculated as a double sum:^{[1]:20}

$$I(X;Y) = \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} p_{(X,Y)}(x,y) \log \left(\frac{p_{(X,Y)}(x,y)}{p_X(x) p_Y(y)} \right), \tag{Eq.1}$$

In terms of PDFs for continuous distributions

In the case of jointly continuous random variables, the double sum is replaced by a double integral:^{[1]:251}

$$I(X;Y) = \int_Y \int_X p_{(X,Y)}(x,y) \log \left(\frac{p_{(X,Y)}(x,y)}{p_X(x)p_Y(y)} \right) dx dy,$$

(Eq.2)

where $p_{(X,Y)}$ is now the joint probability density function of X and Y , and p_X and p_Y are the marginal probability density functions of X and Y respectively.

If the log base 2 is used, the units of mutual information are bits.

Motivation

Intuitively, mutual information measures the information that X and Y share: It measures how much knowing one of these variables reduces uncertainty about the other. For example, if X and Y are independent, then knowing X does not give any information about Y and vice versa, so their mutual information is zero. At the other extreme, if X is a deterministic function of Y and Y is a deterministic function of X then all information conveyed by X is shared with Y : knowing X determines the value of Y and vice versa. As a result, in this case the mutual information is the same as the uncertainty contained in Y (or X) alone, namely the entropy of Y (or X). Moreover, this mutual information is the same as the entropy of X and as the entropy of Y . (A very special case of this is when X and Y are the same random variable.)

Mutual information is a measure of the inherent dependence expressed in the joint distribution of X and Y relative to the joint distribution of X and Y under the assumption of independence. Mutual information therefore measures dependence in the following sense: $I(X;Y) = 0$ if and only if X and Y are independent random variables. This is easy to see in one direction: if X and Y are independent, then $p_{(X,Y)}(x,y) = p_X(x) \cdot p_Y(y)$, and therefore:

$$\log \left(\frac{p_{(X,Y)}(x,y)}{p_X(x)p_Y(y)} \right) = \log 1 = 0.$$

Moreover, mutual information is nonnegative (i.e. $I(X;Y) \geq 0$ see below) and symmetric (i.e. $I(X;Y) = I(Y;X)$ see below).

Relation to other quantities

Nonnegativity

Using Jensen's inequality on the definition of mutual information we can show that $I(X;Y)$ is non-negative, i.e.^{[1]:28}

$$I(X;Y) \geq 0$$

Symmetry

$$I(X;Y) = I(Y;X)$$

Relation to conditional and joint entropy

Mutual information can be equivalently expressed as:

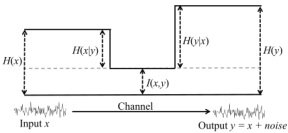
$$\begin{aligned} I(X;Y) &\equiv H(X) - H(X|Y) \\ &\equiv H(Y) - H(Y|X) \\ &\equiv H(X) + H(Y) - H(X,Y) \\ &\equiv H(X,Y) - H(X|Y) - H(Y|X) \end{aligned}$$

where $H(X)$ and $H(Y)$ are the marginal entropies, $H(X|Y)$ and $H(Y|X)$ are the conditional entropies, and $H(X,Y)$ is the joint entropy of X and Y .

Notice the analogy to the union, difference, and intersection of two sets: in this respect, all the formulas given above are apparent from the Venn diagram reported at the beginning of the article.

In terms of a communication channel in which the output Y is a noisy version of the input X , these relations are summarised in the figure:

Because $I(X;Y)$ is non-negative, consequently, $H(X) \geq H(X|Y)$. Here we give the detailed deduction of $I(X;Y) = H(Y) - H(Y|X)$ for the case of jointly discrete random variables:



The relationships between information theoretic quantities

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

The proofs of the other identities above are similar. The proof of the general case (not just discrete) is similar, with integrals replacing sums.

Intuitively, if entropy $H(Y)$ is regarded as a measure of uncertainty about a random variable, then $H(Y|X)$ is a measure of what X does *not* say about Y . This is "the amount of uncertainty remaining about Y after X is known", and thus the right side of the second of these equalities can be read as "the amount of uncertainty in Y , minus the amount of uncertainty in Y which remains after X is known", which is equivalent to "the amount of uncertainty in Y which is removed by knowing X ". This corroborates the intuitive meaning of mutual information as the amount of information (that is, reduction in uncertainty) that knowing either variable provides about the other.

Note that in the discrete case $H(X|X) = 0$ and therefore $H(X) = I(X; X)$. Thus $I(X; X) \geq I(X; Y)$, and one can formulate the basic principle that a variable contains at least as much information about itself as any other variable can provide.

Relation to Kullback–Leibler divergence

For jointly discrete or jointly continuous pairs (X, Y) , mutual information is the Kullback–Leibler divergence of the product of the marginal distributions, $p_X \cdot p_Y$, from the joint distribution $p_{(X,Y)}$, that is,

$$I(X; Y) = D_{\text{KL}}(p_{(X,Y)} \parallel p_X p_Y)$$

Furthermore, let $p_{X|Y=y}(x) = p_{(X,Y)}(x, y)/p_Y(y)$ be the conditional mass or density function. Then, we have the identity

$$I(X; Y) = \mathbb{E}_Y [D_{\text{KL}}(p_{X|Y} \parallel p_X)]$$

The proof for jointly discrete random variables is as follows:

$$\begin{aligned}
I(X; Y) &= \sum_{y \in \mathcal{Y}} p_Y(y) \sum_{x \in \mathcal{X}} p_{X|Y=y}(x) \log \frac{p_{X|Y=y}(x)}{p_X(x)} \\
&= \sum_{y \in \mathcal{Y}} p_Y(y) D_{\text{KL}}(p_{X|Y=y} \parallel p_X) \\
&= \mathbb{E}_Y [D_{\text{KL}}(p_{X|Y} \parallel p_X)].
\end{aligned}$$

Similarly this identity can be established for jointly continuous random variables.

Note that here the Kullback–Leibler divergence involves integration over the values of the random variable X only, and the expression $D_{\text{KL}}(p_{X|Y} \parallel p_X)$ still denotes a random variable because Y is random. Thus mutual information can also be understood as the expectation of the Kullback–Leibler divergence of the univariate distribution p_X of X from the conditional distribution $p_{X|Y}$ of X given Y : the more different the distributions $p_{X|Y}$ and p_X are on average, the greater the information gain.

Bayesian estimation of mutual information

It is well-understood how to do Bayesian estimation of the mutual information of a joint distribution based on samples of that distribution. The first work to do this, which also showed how to do Bayesian estimation of many other information-theoretic properties besides mutual information, was [2]. Subsequent researchers have rederived [3] and extended [4] this analysis. See [5] for a recent paper based on a prior specifically tailored to estimation of mutual information per se. Besides, recently an estimation method accounting for continuous and multivariate outputs, Y , was proposed in [6].

Independence assumptions

The Kullback–Leibler divergence formulation of the mutual information is predicated on that one is interested in comparing $p(x, y)$ to the fully factorized outer product $p(x) \cdot p(y)$. In many problems, such as non-negative matrix factorization, one is interested in less extreme factorizations; specifically, one wishes to compare $p(x, y)$ to a low-rank matrix approximation in some unknown variable w ; that is, to what degree one might have

$$p(x, y) \approx \sum_w p'(x, w) p''(w, y)$$

Alternately, one might be interested in knowing how much more information $\mathbf{p}(\mathbf{x}, \mathbf{y})$ carries over its factorization. In such a case, the excess information that the full distribution $\mathbf{p}(\mathbf{x}, \mathbf{y})$ carries over the matrix factorization is given by the Kullback-Leibler divergence

$$\mathbf{I}_{LRMA} = \sum_{\mathbf{y} \in \mathcal{Y}} \sum_{\mathbf{x} \in \mathcal{X}} \mathbf{p}(\mathbf{x}, \mathbf{y}) \log \left(\frac{\mathbf{p}(\mathbf{x}, \mathbf{y})}{\sum_{\mathbf{w}} \mathbf{p}'(\mathbf{x}, \mathbf{w}) \mathbf{p}''(\mathbf{w}, \mathbf{y})} \right),$$

The conventional definition of the mutual information is recovered in the extreme case that the process \mathbf{W} has only one value for \mathbf{w} .

Variations

Several variations on mutual information have been proposed to suit various needs. Among these are normalized variants and generalizations to more than two variables.

Metric

Many applications require a metric, that is, a distance measure between pairs of points. The quantity

$$\begin{aligned} d(\mathbf{X}, \mathbf{Y}) &= \mathbf{H}(\mathbf{X}, \mathbf{Y}) - \mathbf{I}(\mathbf{X}; \mathbf{Y}) \\ &= \mathbf{H}(\mathbf{X}) + \mathbf{H}(\mathbf{Y}) - 2\mathbf{I}(\mathbf{X}; \mathbf{Y}) \\ &= \mathbf{H}(\mathbf{X}|\mathbf{Y}) + \mathbf{H}(\mathbf{Y}|\mathbf{X}) \end{aligned}$$

satisfies the properties of a metric (triangle inequality, non-negativity, indiscernability and symmetry). This distance metric is also known as the variation of information.

If \mathbf{X}, \mathbf{Y} are discrete random variables then all the entropy terms are non-negative, so $0 \leq d(\mathbf{X}, \mathbf{Y}) \leq \mathbf{H}(\mathbf{X}, \mathbf{Y})$ and one can define a normalized distance

$$D(\mathbf{X}, \mathbf{Y}) = \frac{d(\mathbf{X}, \mathbf{Y})}{\mathbf{H}(\mathbf{X}, \mathbf{Y})} \leq 1.$$

The metric D is a universal metric, in that if any other distance measure places \mathbf{X} and \mathbf{Y} close-by, then the D will also judge them close.^[7]

Plugging in the definitions shows that

$$D(\mathbf{X}, \mathbf{Y}) = 1 - \frac{\mathbf{I}(\mathbf{X}; \mathbf{Y})}{\mathbf{H}(\mathbf{X}, \mathbf{Y})}.$$

In a set-theoretic interpretation of information (see the figure for Conditional entropy), this is effectively the Jaccard distance between \mathbf{X} and \mathbf{Y} .

Finally,

$$D'(\mathbf{X}, \mathbf{Y}) = 1 - \frac{\mathbf{I}(\mathbf{X}; \mathbf{Y})}{\max\{\mathbf{H}(\mathbf{X}), \mathbf{H}(\mathbf{Y})\}}$$

is also a metric.

Conditional mutual information

Sometimes it is useful to express the mutual information of two random variables conditioned on a third.

$$\mathbf{I}(\mathbf{X}; \mathbf{Y}|\mathbf{Z}) = \mathbb{E}_{\mathbf{Z}}[D_{\text{KL}}(P_{(\mathbf{X}, \mathbf{Y})|\mathbf{Z}} \| P_{\mathbf{X}|\mathbf{Z}} \otimes P_{\mathbf{Y}|\mathbf{Z}})]$$

For jointly discrete random variables this takes the form

$$\mathbf{I}(\mathbf{X}; \mathbf{Y}|\mathbf{Z}) = \sum_{\mathbf{z} \in \mathcal{Z}} \sum_{\mathbf{y} \in \mathcal{Y}} \sum_{\mathbf{x} \in \mathcal{X}} p_{\mathbf{Z}}(\mathbf{z}) p_{\mathbf{X}, \mathbf{Y}|\mathbf{Z}}(\mathbf{x}, \mathbf{y}|\mathbf{z}) \log \left[\frac{p_{\mathbf{X}, \mathbf{Y}|\mathbf{Z}}(\mathbf{x}, \mathbf{y}|\mathbf{z})}{p_{\mathbf{X}|\mathbf{Z}}(\mathbf{x}|\mathbf{z}) p_{\mathbf{Y}|\mathbf{Z}}(\mathbf{y}|\mathbf{z})} \right],$$

which can be simplified as

$$\mathbf{I}(\mathbf{X}; \mathbf{Y}|\mathbf{Z}) = \sum_{\mathbf{z} \in \mathcal{Z}} \sum_{\mathbf{y} \in \mathcal{Y}} \sum_{\mathbf{x} \in \mathcal{X}} p_{\mathbf{X}, \mathbf{Y}, \mathbf{Z}}(\mathbf{x}, \mathbf{y}, \mathbf{z}) \log \frac{p_{\mathbf{X}, \mathbf{Y}, \mathbf{Z}}(\mathbf{x}, \mathbf{y}, \mathbf{z}) p_{\mathbf{Z}}(\mathbf{z})}{p_{\mathbf{X}, \mathbf{Z}}(\mathbf{x}, \mathbf{z}) p_{\mathbf{Y}, \mathbf{Z}}(\mathbf{y}, \mathbf{z})}.$$

For jointly continuous random variables this takes the form

$$\mathbf{I}(\mathbf{X}; \mathbf{Y}|\mathbf{Z}) = \int_{\mathcal{Z}} \int_{\mathcal{Y}} \int_{\mathcal{X}} p_{\mathbf{Z}}(\mathbf{z}) p_{\mathbf{X}, \mathbf{Y}|\mathbf{Z}}(\mathbf{x}, \mathbf{y}|\mathbf{z}) \log \left[\frac{p_{\mathbf{X}, \mathbf{Y}|\mathbf{Z}}(\mathbf{x}, \mathbf{y}|\mathbf{z})}{p_{\mathbf{X}|\mathbf{Z}}(\mathbf{x}|\mathbf{z}) p_{\mathbf{Y}|\mathbf{Z}}(\mathbf{y}|\mathbf{z})} \right] d\mathbf{x} d\mathbf{y} d\mathbf{z},$$

which can be simplified as

$$\mathbf{I}(\mathbf{X}; \mathbf{Y}|\mathbf{Z}) = \int_{\mathcal{Z}} \int_{\mathcal{Y}} \int_{\mathcal{X}} p_{\mathbf{X}, \mathbf{Y}, \mathbf{Z}}(\mathbf{x}, \mathbf{y}, \mathbf{z}) \log \frac{p_{\mathbf{X}, \mathbf{Y}, \mathbf{Z}}(\mathbf{x}, \mathbf{y}, \mathbf{z}) p_{\mathbf{Z}}(\mathbf{z})}{p_{\mathbf{X}, \mathbf{Z}}(\mathbf{x}, \mathbf{z}) p_{\mathbf{Y}, \mathbf{Z}}(\mathbf{y}, \mathbf{z})} d\mathbf{x} d\mathbf{y} d\mathbf{z}.$$

Conditioning on a third random variable may either increase or decrease the mutual information, but it is always true that

$$\mathbf{I}(\mathbf{X}; \mathbf{Y}|\mathbf{Z}) \geq 0$$

for discrete, jointly distributed random variables $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$. This result has been used as a basic building block for proving other inequalities in information theory.

Multivariate mutual information

Several generalizations of mutual information to more than two random variables have been proposed, such as total correlation (or multi-information) and interaction information. The expression and study of multivariate higher-degree mutual-information was achieved in two seemingly independent works: McGill (1954)^[8] who called these functions “interaction information”, and Hu Kuo Ting (1962)^[9] who also first proved the possible negativity of mutual-information for degrees higher than 2 and justified algebraically the intuitive correspondence to Venn diagrams^[10]

$$I(\mathbf{X}_1; \mathbf{X}_1) = H(\mathbf{X}_1)$$

and for $n > 1$,

$$I(\mathbf{X}_1; \dots; \mathbf{X}_n) = I(\mathbf{X}_1; \dots; \mathbf{X}_{n-1}) - I(\mathbf{X}_1; \dots; \mathbf{X}_{n-1} | \mathbf{X}_n),$$

where (as above) we define

$$I(\mathbf{X}_1; \dots; \mathbf{X}_{n-1} | \mathbf{X}_n) = \mathbb{E}_{\mathbf{X}_n} [D_{\text{KL}}(P_{(\mathbf{X}_1, \dots, \mathbf{X}_{n-1}) | \mathbf{X}_n} \| P_{\mathbf{X}_1 | \mathbf{X}_n} \otimes \dots \otimes P_{\mathbf{X}_{n-1} | \mathbf{X}_n})].$$

(This definition of multivariate mutual information is identical to that of interaction information except for a change in sign when the number of random variables is odd.)

Multivariate statistical independence

The multivariate mutual-information functions generalize the pairwise independence case that states that $\mathbf{X}_1, \mathbf{X}_2$ if and only if $I(\mathbf{X}_1; \mathbf{X}_2) = 0$, to arbitrary numerous variable. n variables are mutually independent if and only if the $2^n - n - 1$ mutual information functions vanish $I(\mathbf{X}_1; \dots; \mathbf{X}_k) = 0$ with $n \geq k \geq 2$ (theorem 2^[10]). In this sense, the $I(\mathbf{X}_1; \dots; \mathbf{X}_k) = 0$ can be used as a refined statistical independence criterion.

Applications

For 3 variables, Brenner et al. applied multivariate mutual information to neural coding and called its negativity "synergy"^[11] and Watkinson et al. applied it to genetic expression^[12]. For arbitrary k variables, Tapia et al. applied multivariate mutual information to gene expression^[13]^[10]. It can be zero, positive, or negative^[14]. The positivity corresponds to relations generalizing the pairwise correlations, nullity corresponds to a refined notion of independence, and negativity detects high dimensional "emergent" relations and clusterized datapoints^[13].

One high-dimensional generalization scheme which maximizes the mutual information between the joint distribution and other target variables is found to be useful in feature selection.^[15]

Mutual information is also used in the area of signal processing as a measure of similarity between two signals. For example, FMI metric^[16] is an image fusion performance measure that makes use of mutual information in order to measure the amount of information that the fused image contains about the source images. The Matlab code for this metric can be found at.^[17]

Directed information

Directed information, $I(\mathbf{X}^n \rightarrow \mathbf{Y}^n)$, measures the amount of information that flows from the process \mathbf{X}^n to \mathbf{Y}^n , where \mathbf{X}^n denotes the vector $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ and \mathbf{Y}^n denotes $\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_n$. The term *directed information* was coined by James Massey and is defined as

$$I(\mathbf{X}^n \rightarrow \mathbf{Y}^n) = \sum_{i=1}^n I(\mathbf{X}^i; \mathbf{Y}_i | \mathbf{Y}^{i-1}).$$

Note that if $n = 1$, the directed information becomes the mutual information. Directed information has many applications in problems where causality plays an important role, such as capacity of channel with feedback.^{[18][19]}

Normalized variants

Normalized variants of the mutual information are provided by the *coefficients of constraint*,^[20] uncertainty coefficient^[21] or proficiency:^[22]

$$C_{XY} = \frac{I(\mathbf{X}; \mathbf{Y})}{H(\mathbf{Y})} \quad \text{and} \quad C_{YX} = \frac{I(\mathbf{X}; \mathbf{Y})}{H(\mathbf{X})}.$$

The two coefficients have a value ranging in [0, 1], but are not necessarily equal. In some cases a symmetric measure may be desired, such as the following redundancy measure:

$$R = \frac{I(\mathbf{X}; \mathbf{Y})}{H(\mathbf{X}) + H(\mathbf{Y})}$$

which attains a minimum of zero when the variables are independent and a maximum value of

$$R_{\max} = \frac{\min \{H(\mathbf{X}), H(\mathbf{Y})\}}{H(\mathbf{X}) + H(\mathbf{Y})}$$

when one variable becomes completely redundant with the knowledge of the other. See also *Redundancy (information theory)*.

Another symmetrical measure is the *symmetric uncertainty* (Witten & Frank 2005), given by

$$U(X, Y) = 2R = 2 \frac{I(X; Y)}{H(X) + H(Y)}$$

which represents the harmonic mean of the two uncertainty coefficients C_{XY}, C_{YX} .^[21]

If we consider mutual information as a special case of the total correlation or dual total correlation, the normalized version are respectively,

$$\frac{I(X; Y)}{\min [H(X), H(Y)]} \text{ and } \frac{I(X; Y)}{H(X, Y)}.$$

This normalized version also known as **Information Quality Ratio (IQR)** which quantifies the amount of information of a variable based on another variable against total uncertainty.^[23]

$$IQR(X, Y) = E[I(X; Y)] = \frac{I(X; Y)}{H(X, Y)} = \frac{\sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x) p(y)}{\sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x, y)} - 1$$

There's a normalization^[24] which derives from first thinking of mutual information as an analogue to covariance (thus Shannon entropy is analogous to variance). Then the normalized mutual information is calculated akin to the Pearson correlation coefficient,

$$\frac{I(X; Y)}{\sqrt{H(X)H(Y)}}.$$

Weighted variants

In the traditional formulation of the mutual information,

$$I(X; Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \frac{p(x, y)}{p(x) p(y)},$$

each *event* or *object* specified by (x, y) is weighted by the corresponding probability $p(x, y)$. This assumes that all objects or events are equivalent *apart from* their probability of occurrence. However, in some applications it may be the case that certain objects or events are more *significant* than others, or that certain patterns of association are more semantically important than others.

For example, the deterministic mapping $\{(1, 1), (2, 2), (3, 3)\}$ may be viewed as stronger than the deterministic mapping $\{(1, 3), (2, 1), (3, 2)\}$, although these relationships would yield the same mutual information. This is because the mutual information is not sensitive at all to any inherent ordering in the variable values (Cronbach 1954, Coombs, Dawes & Tversky 1970, Lockhead 1970), and is therefore not sensitive at all to the **form** of the relational mapping between the associated variables. If it is desired that the former relation—showing agreement on all variable values—be judged stronger than the later relation, then it is possible to use the following *weighted mutual information* (Guíasu 1977).

$$I(X; Y) = \sum_{y \in Y} \sum_{x \in X} w(x, y) p(x, y) \log \frac{p(x, y)}{p(x) p(y)},$$

which places a weight $w(x, y)$ on the probability of each variable value co-occurrence, $p(x, y)$. This allows that certain probabilities may carry more or less significance than others, thereby allowing the quantification of relevant *holistic* or *Prägnanz* factors. In the above example, using larger relative weights for $w(1, 1)$, $w(2, 2)$, and $w(3, 3)$ would have the effect of assessing greater *informativeness* for the relation $\{(1, 1), (2, 2), (3, 3)\}$ than for the relation $\{(1, 3), (2, 1), (3, 2)\}$, which may be desirable in some cases of pattern recognition, and the like. This weighted mutual information is a form of weighted KL-Divergence, which is known to take negative values for some inputs,^[25] and there are examples where the weighted mutual information also takes negative values.^[26]

Adjusted mutual information

A probability distribution can be viewed as a partition of a set. One may then ask: if a set were partitioned randomly, what would the distribution of probabilities be? What would the expectation value of the mutual information be? The adjusted mutual information or AMI subtracts the expectation value of the MI, so that the AMI is zero when two different distributions are random, and one when two distributions are identical. The AMI is defined in analogy to the adjusted Rand index of two different partitions of a set.

Absolute mutual information

Using the ideas of Kolmogorov complexity, one can consider the mutual information of two sequences independent of any probability distribution:

$$I_K(X; Y) = K(X) - K(X|Y).$$

To establish that this quantity is symmetric up to a logarithmic factor ($I_K(X; Y) \approx I_K(Y; X)$) one requires the chain rule for Kolmogorov complexity (Li & Vitányi 1997). Approximations of this quantity via compression can be used to define a distance measure to perform a hierarchical clustering of sequences without having any domain knowledge of the sequences (Cilibiasi & Vitányi 2005).

Linear correlation

Unlike correlation coefficients, such as the product moment correlation coefficient, mutual information contains information about all dependence—linear and nonlinear—and not just linear dependence as the correlation coefficient measures. However, in the narrow case that the joint distribution for X and Y is a bivariate normal distribution (implying in particular that both marginal distributions are normally distributed), there is an exact relationship between I and the correlation coefficient ρ (Gel'fand & Yaglom 1957).

$$I = -\frac{1}{2} \log(1 - \rho^2)$$

The equation above can be derived as follows for a bivariate Gaussian:

$$\begin{aligned} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} &\sim \mathcal{N}\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \Sigma\right), \quad \Sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix} \\ \mathbf{H}(X_i) &= \frac{1}{2} \log(2\pi e \sigma_i^2) = \frac{1}{2} + \frac{1}{2} \log(2\pi) + \log(\sigma_i), \quad i \in \{1, 2\} \\ \mathbf{H}(X_1, X_2) &= \frac{1}{2} \log[(2\pi e)^2 |\Sigma|] = 1 + \log(2\pi) + \log(\sigma_1\sigma_2) + \frac{1}{2} \log(1 - \rho^2) \end{aligned}$$

Therefore,

$$\mathbf{I}(X_1; X_2) = \mathbf{H}(X_1) + \mathbf{H}(X_2) - \mathbf{H}(X_1, X_2) = -\frac{1}{2} \log(1 - \rho^2)$$

For discrete data

When **X** and **Y** are limited to be in a discrete number of states, observation data is summarized in a contingency table, with row variable **X** (or *i*) and column variable **Y** (or *j*). Mutual information is one of the measures of association or correlation between the row and column variables. Other measures of association include Pearson's chi-squared test statistics, G-test statistics, etc. In fact, mutual information is equal to G-test statistics divided by **2N**, where **N** is the sample size.

Applications

In many applications, one wants to maximize mutual information (thus increasing dependencies), which is often equivalent to minimizing conditional entropy. Examples include:

- In search engine technology, mutual information between phrases and contexts is used as a feature for k-means clustering to discover semantic clusters (concepts).^[27] For example, the mutual information of a bigram might be calculated as:

$$MI(x, y) = \log \frac{P_{X,Y}(x, y)}{P_X(x)P_Y(y)} \approx \log \frac{\frac{f_{xy}}{B}}{\frac{f_x}{U} \frac{f_y}{U}}$$

where *f_{xy}* is the number of times the bigram xy appears in the corpus, *f_x* is the number of times the unigram x appears in the corpus, B is the total number of bigrams, and U is the total number of unigrams.^[27]

- In telecommunications, the channel capacity is equal to the mutual information, maximized over all input distributions.
- Discriminative training procedures for hidden Markov models have been proposed based on the maximum mutual information (MMI) criterion.
- RNA secondary structure prediction from a multiple sequence alignment.
- Phylogenetic profiling prediction from pairwise present and disappearance of functionally link genes.
- Mutual information has been used as a criterion for feature selection and feature transformations in machine learning. It can be used to characterize both the relevance and redundancy of variables, such as the minimum redundancy feature selection.
- Mutual information is used in determining the similarity of two different clusterings of a dataset. As such, it provides some advantages over the traditional Rand index.
- Mutual information of words is often used as a significance function for the computation of collocations in corpus linguistics. This has the added complexity that no word-instance is an instance to two different words; rather, one counts instances where 2 words occur adjacent or in close proximity; this slightly complicates the calculation, since the expected probability of one word occurring within **N** words of another, goes up with **N**.
- Mutual information is used in medical imaging for image registration. Given a reference image (for example, a brain scan), and a second image which needs to be put into the same coordinate system as the reference image, this image is deformed until the mutual information between it and the reference image is maximized.
- Detection of phase synchronization in time series analysis
- In the infomax method for neural-net and other machine learning, including the infomax-based Independent component analysis algorithm
- Average mutual information in delay embedding theorem is used for determining the *embedding delay* parameter.
- Mutual information between genes in expression microarray data is used by the ARACNE algorithm for reconstruction of gene networks.
- In statistical mechanics, Loschmidt's paradox may be expressed in terms of mutual information.^{[28][29]} Loschmidt noted that it must be impossible to determine a physical law which lacks time reversal symmetry (e.g. the second law of thermodynamics) only from physical laws which have this symmetry. He pointed out that the H-theorem of Boltzmann made the assumption that the velocities of particles in a gas were permanently uncorrelated, which removed the time symmetry inherent in the H-theorem. It can be shown that if a system is described by a probability density in phase space, then Liouville's theorem implies that the joint information (negative of the joint entropy) of the distribution remains constant in time. The joint information is equal to the mutual information plus the sum of all the marginal information (negative of the marginal entropies) for each particle coordinate. Boltzmann's assumption amounts to ignoring the mutual information in the calculation of entropy, which yields the thermodynamic entropy (divided by Boltzmann's constant).
- The mutual information is used to learn the structure of Bayesian networks/dynamic Bayesian networks, which is thought to explain the causal relationship between random variables, as exemplified by the GlobalMIT toolkit:^[30] learning the globally optimal dynamic Bayesian network with the Mutual Information Test criterion.
- Popular cost function in decision tree learning.
- The mutual information is used in cosmology to test the influence of large-scale environments on galaxy properties in the Galaxy Zoo.
- The mutual information was used in Solar Physics to derive the solar differential rotation profile, a travel-time deviation map for sunspots, and a time–distance diagram from quiet-Sun measurements^[31]
- Used in Invariant Information Clustering to automatically train neural network classifiers and image segmenters given no labelled data.^[32]

See also

- Pointwise mutual information
- Quantum mutual information

Notes

1. Cover, T.M.; Thomas, J.A. (1991). *Elements of Information Theory* (<http://archive.org/details/elementsofinform0000cove>) (Wiley ed.). ISBN 978-0-471-24195-9.
2. Wolpert, D.H.; Wolf, D.R. (1995). "Estimating functions of probability distributions from a finite set of samples". *Physical Review E*. **52** (6): 6841–6854. Bibcode:1995PhRvE..52.6841W (<https://ui.adsabs.harvard.edu/abs/1995PhRvE..52.6841W>). CiteSeerX 10.1.1.55.7122 (<https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.55.7122>). doi:10.1103/PhysRevE.52.6841 (<https://doi.org/10.1103%2FPhysRevE.52.6841>). PMID 9964199 (<https://pubmed.ncbi.nlm.nih.gov/9964199/>).
3. Hutter, M. (2001). "Distribution of Mutual Information". *Advances in Neural Information Processing Systems 2001*.
4. Archer, E.; Park, I.M.; Pillow, J. (2013). "Bayesian and Quasi-Bayesian Estimators for Mutual Information from Discrete Data". *Entropy*. **15** (12): 1738–1755. Bibcode:2013Entpr..15.1738A (<https://ui.adsabs.harvard.edu/abs/2013Entpr..15.1738A>). CiteSeerX 10.1.1.294.4690 (<https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.294.4690>). doi:10.3390/e15051738 (<https://doi.org/10.3390%2Fe15051738>).
5. Wolpert, D.H.; DeDeo, S. (2013). "Estimating Functions of Distributions Defined over Spaces of Unknown Size". *Entropy*. **15** (12): 4668–4699. arXiv:1311.4548 (<https://arxiv.org/abs/1311.4548>). Bibcode:2013Entpr..15.4668W (<https://ui.adsabs.harvard.edu/abs/2013Entpr..15.4668W>). doi:10.3390/e15114668 (<https://doi.org/10.3390%2Fe15114668>).
6. Tomasz Jetka; Karol Nienaltowski; Tomasz Winarski; Sławomir Blonski; Michał Komorowski (2019), "Information-theoretic analysis of multivariate single-cell signaling responses", *PLOS Computational Biology*, **15** (7): e1007132, arXiv:1808.05581 (<https://arxiv.org/abs/1808.05581>), Bibcode:2019PLSCB..15E7132J (<https://ui.adsabs.harvard.edu/abs/2019PLSCB..15E7132J>), doi:10.1371/journal.pcbi.1007132 (<https://doi.org/10.1371%2Fjournal.pcbi.1007132>), PMC 6655862 (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6655862>), PMID 31299056 (<https://pubmed.ncbi.nlm.nih.gov/31299056/>)
7. Kraskov, Alexander; Stögbauer, Harald; Andrzejak, Ralph G.; Grassberger, Peter (2003). "Hierarchical Clustering Based on Mutual Information". arXiv:q-bio/0311039 (<https://arxiv.org/abs/q-bio/0311039>). Bibcode:2003q.bio....11039K (<https://ui.adsabs.harvard.edu/abs/2003q.bio....11039K>).
8. McGill, W. (1954). "Multivariate information transmission". *Psychometrika*. **19** (1): 97–116. doi:10.1007/BF02289159 (<https://doi.org/10.1007%2FBF02289159>).
9. Hu, K.T. (1962). "On the Amount of Information". *Theory Probab. Appl.* **7** (4): 439–447. doi:10.1137/1107041 (<https://doi.org/10.1137%2F1107041>).
10. Baudot, P.; Tapia, M.; Bennequin, D.; Goillard, J.M. (2019). "Topological Information Data Analysis". *Entropy*. **21** (9): 869. arXiv:1907.04242 (<https://arxiv.org/abs/1907.04242>). Bibcode:2019Entpr..21..869B (<https://ui.adsabs.harvard.edu/abs/2019Entpr..21..869B>). doi:10.3390/e21090869 (<https://doi.org/10.3390%2Fe21090869>).
11. Brenner, N.; Strong, S.; Koberle, R.; Bialek, W. (2000). "Synergy in a Neural Code". *Neural Comput.* **12** (7): 1531–1552. doi:10.1162/0899766000300015259 (<https://doi.org/10.1162%2F089976600300015259>). PMID 10935917 (<https://pubmed.ncbi.nlm.nih.gov/10935917/>).
12. Watkinson, J.; Liang, K.; Wang, X.; Zheng, T.; Anastassiou, D. (2009). "Inference of Regulatory Gene Interactions from Expression Data Using Three-Way Mutual Information" (<https://semanticscholar.org/paper/cb09223a34b08e6dcbf696385d9ab76fd9f37aa4>). *Chall. Syst. Biol. Ann. N. Y. Acad. Sci.* **1158** (1): 302–313. Bibcode:2009NYASA1158..302W (<https://ui.adsabs.harvard.edu/abs/2009NYASA1158..302W>). doi:10.1111/j.1749-6632.2008.03757.x (<https://doi.org/10.1111%2Fj.1749-6632.2008.03757.x>). PMID 19348651 (<https://pubmed.ncbi.nlm.nih.gov/19348651/>).
13. Tapia, M.; Baudot, P.; Formizano-Treziny, C.; Dufour, M.; Goillard, J.M. (2018). "Neurotransmitter identity and electrophysiological phenotype are genetically coupled in midbrain dopaminergic neurons" (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6134142>). *Sci. Rep.* **8** (1): 13637. Bibcode:2018NatSR...813637T (<https://ui.adsabs.harvard.edu/abs/2018NatSR...813637T>). doi:10.1038/s41598-018-31765-z (<https://doi.org/10.1038%2Fs41598-018-31765-z>). PMC 6134142 (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6134142>). PMID 30206240 (<https://pubmed.ncbi.nlm.nih.gov/30206240/>).
14. Hu, K.T. (1962). "On the Amount of Information". *Theory Probab. Appl.* **7** (4): 439–447. doi:10.1137/1107041 (<https://doi.org/10.1137%2F1107041>).
15. Christopher D. Manning; Prabhakar Raghavan; Hinrich Schütze (2008). *An Introduction to Information Retrieval*. Cambridge University Press. ISBN 978-0-521-86571-5.
16. Haghighat, M. B. A.; Aghagolzadeh, A.; Seyedarabi, H. (2011). "A non-reference image fusion metric based on mutual information of image features". *Computers & Electrical Engineering*. **37** (5): 744–756. doi:10.1016/j.compeleceng.2011.07.012 (<https://doi.org/10.1016%2Fj.compeleceng.2011.07.012>).
17. "Feature Mutual Information (FMI) metric for non-reference image fusion - File Exchange - MATLAB Central" (<http://www.mathworks.com/matlabcentral/fileexchange/45926-feature-mutual-information-fmi-image-fusion-metric>). *www.mathworks.com*. Retrieved 4 April 2018.
18. Massey, James (1990). "Causality, Feedback And Directed Information". *Proc. 1990 Intl. Symp. on Info. Th. and its Applications, Waikiki, Hawaii, Nov. 27-30, 1990*. CiteSeerX 10.1.1.36.5688 (<https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.36.5688>).
19. Permuter, Haim Henry; Weissman, Tsachy; Goldsmith, Andrea J. (February 2009). "Finite State Channels With Time-Invariant Deterministic Feedback". *IEEE Transactions on Information Theory*. **55** (2): 644–662. arXiv:cs/0608070 (<https://arxiv.org/abs/cs/0608070>). doi:10.1109/TIT.2008.2009849 (<https://doi.org/10.1109%2FTIT.2008.2009849>).
20. Coombs, Dawes & Tversky 1970.
21. Press, WH; Teukolsky, SA; Vetterling, WT; Flannery, BP (2007). "Section 14.7.3. Conditional Entropy and Mutual Information" (<http://ap.ps.nrbook.com/empanel/index.html#pg=758>). *Numerical Recipes: The Art of Scientific Computing* (3rd ed.). New York: Cambridge University Press. ISBN 978-0-521-88068-8.
22. White, Jim; Steingold, Sam; Fournelle, Connie. *Performance Metrics for Group-Detection Algorithms* (<http://www.interfacesymposia.org/I04/I2004Proceedings/WhiteJim/WhiteJim.paper.pdf>) (PDF). Interface 2004.
23. Wijaya, Dedy Rahman; Sarno, Riyanarto; Zulaika, Enny (2017). "Information Quality Ratio as a novel metric for mother wavelet selection". *Chemometrics and Intelligent Laboratory Systems*. **160**: 59–71. doi:10.1016/j.chemolab.2016.11.012 (<https://doi.org/10.1016%2Fj.chemolab.2016.11.012>).
24. Strehl, Alexander; Ghosh, Joydeep (2003). "Cluster Ensembles – A Knowledge Reuse Framework for Combining Multiple Partitions" (<http://www.jmlr.org/papers/volume3/strehl02a/strehl02a.pdf>) (PDF). *The Journal of Machine Learning Research*. **3**: 583–617. doi:10.1162/153244303321897735 (<https://doi.org/10.1162%2F153244303321897735>).
25. Kvålseth, T. O. (1991). "The relative useful information measure: some comments". *Information Sciences*. **56** (1): 35–38. doi:10.1016/0020-0255(91)90022-m (<https://doi.org/10.1016%2F0020-0255%2891%2990022-m>).
26. Pocock, A. (2012). *Feature Selection Via Joint Likelihood* (<http://www.cs.man.ac.uk/~gbrown/publications/pocockPhDthesis.pdf>) (PDF) (Thesis).
27. Parsing a Natural Language Using Mutual Information Statistics (<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.78.4178&rep=rep1&type=pdf>) by David M. Magerman and Mitchell P. Marcus
28. Hugh Everett Theory of the Universal Wavefunction (<https://www.pbs.org/wgbh/nova/manyworlds/pdf/dissertation.pdf>), Thesis, Princeton University, (1956, 1973), pp 1–140 (page 30)
29. Everett, Hugh (1957). "Relative State Formulation of Quantum Mechanics" (<https://web.archive.org/web/20111027191052/http://www.univer.omsk.su/omsk/Sci/Everett/paper1957.html>). *Reviews of Modern Physics*. **29** (3): 454–462. Bibcode:1957RvMP...29..454E (<https://ui.adsabs.harvard.edu/abs/1957RvMP...29..454E>). doi:10.1103/revmodphys.29.454 (<https://doi.org/10.1103%2Frevmodphys.29.454>). Archived from the original (<http://www.univer.omsk.su/omsk/Sci/Everett/paper1957.html>) on 2011-10-27. Retrieved 2012-07-16.
30. GlobalMIT (<https://code.google.com/p/globalmit>) at Google Code
31. Keys, Dustin; Kholikov, Shukur; Pevtsov, Alexei A. (February 2015). "Application of Mutual Information Methods in Time Distance Helioseismology". *Solar Physics*. **290** (3): 659–671. arXiv:1501.05597 (<https://arxiv.org/abs/1501.05597>). Bibcode:2015SoPh..290..659K (<https://ui.adsabs.harvard.edu/abs/2015SoPh..290..659K>). doi:10.1007/s11207-015-0650-y (<https://doi.org/10.1007%2Fs11207-015-0650-y>).
32. Invariant Information Clustering for Unsupervised Image Classification and Segmentation (<https://arxiv.org/abs/1807.06653>) by Xu Ji, Joao Henriques and Andrea Vedaldi

References

- Baudot, P.; Tapia, M.; Bennequin, D.; Goiaillard, J.M. (2019). "Topological Information Data Analysis". *Entropy*. **21** (9). 869. arXiv:1907.04242 (<https://arxiv.org/abs/1907.04242>). Bibcode:2019Entrp..21..869B (<https://ui.adsabs.harvard.edu/abs/2019Entrp..21..869B>). doi:10.3390/e21090869 (<https://doi.org/10.3390/e21090869>).
- Cilibrasi, R.; Vitányi, Paul (2005). "Clustering by compression" (<http://www.cwi.nl/~paulv/papers/cluster.pdf>) (PDF). *IEEE Transactions on Information Theory*. **51** (4): 1523–1545. arXiv:cs/0312044 (<https://arxiv.org/abs/cs/0312044>). doi:10.1109/TIT.2005.844059 (<https://doi.org/10.1109%2FTIT.2005.844059>).
- Cronbach, L. J. (1954). "On the non-rational application of information measures in psychology". In Quastler, Henry (ed.). *Information Theory in Psychology: Problems and Methods*. Glencoe, Illinois: Free Press. pp. 14–30.
- Coombs, C. H.; Dawes, R. M.; Tversky, A. (1970). *Mathematical Psychology: An Elementary Introduction*. Englewood Cliffs, New Jersey: Prentice-Hall.
- Church, Kenneth Ward; Hanks, Patrick (1989). "Word association norms, mutual information, and lexicography" (<http://dl.acm.org/citation.cfm?id=89095>). *Proceedings of the 27th Annual Meeting of the Association for Computational Linguistics*: 76–83. doi:10.3115/981623.981633 (<https://doi.org/10.3115%2F981623.981633>).
- Gel'fand, I.M.; Yaglom, A.M. (1957). "Calculation of amount of information about a random function contained in another such function". *American Mathematical Society Translations: Series 2*. **12**: 199–246. English translation of original in *Uspekhi Matematicheskikh Nauk* **12** (1): 3–52.
- Guiasu, Silviu (1977). *Information Theory with Applications*. McGraw-Hill, New York. ISBN 978-0-07-025109-0.
- Li, Ming; Vitányi, Paul (February 1997). *An introduction to Kolmogorov complexity and its applications*. New York: Springer-Verlag. ISBN 978-0-387-94868-3.
- Lockhead, G. R. (1970). "Identification and the form of multidimensional discrimination space". *Journal of Experimental Psychology*. **85** (1): 1–10. doi:10.1037/h0029508 (<https://doi.org/10.1037%2Fh0029508>). PMID 5458322 (<https://pubmed.ncbi.nlm.nih.gov/5458322>).
- David J. C. MacKay. *Information Theory, Inference, and Learning Algorithms* (<http://www.inference.phy.cam.ac.uk/mackay/itila/book.html>) Cambridge: Cambridge University Press, 2003. ISBN 0-521-64298-1 (available free online)
- Haghighat, M. B. A.; Aghagolzadeh, A.; Seyedarabi, H. (2011). "A non-reference image fusion metric based on mutual information of image features". *Computers & Electrical Engineering*. **37** (5): 744–756. doi:10.1016/j.compeleceng.2011.07.012 (<https://doi.org/10.1016%2Fj.compeleceng.2011.07.012>).
- Athanasios Papoulis. *Probability, Random Variables, and Stochastic Processes*, second edition. New York: McGraw-Hill, 1984. (See Chapter 15.)
- Witten, Ian H. & Frank, Eibe (2005). *Data Mining: Practical Machine Learning Tools and Techniques* (<http://www.cs.waikato.ac.nz/~ml/weka/book.html>). Morgan Kaufmann, Amsterdam. ISBN 978-0-12-374856-0.
- Peng, H.C.; Long, F. & Ding, C. (2005). "Feature selection based on mutual information: criteria of max-dependency, max-relevance, and min-redundancy" (<http://research.janelia.org/peng/proj/mRMR/index.htm>). *IEEE Transactions on Pattern Analysis and Machine Intelligence*. **27** (8): 1226–1238. CiteSeerX 10.1.1.63.5765 (<https://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.63.5765>). doi:10.1109/tpami.2005.159 (<https://doi.org/10.1109%2Ftpami.2005.159>). PMID 16119262 (<https://pubmed.ncbi.nlm.nih.gov/16119262>).
- Andre S. Ribeiro; Stuart A. Kauffman; Jason Lloyd-Price; Bjorn Samuelsson & Joshua Socolar (2008). "Mutual Information in Random Boolean models of regulatory networks". *Physical Review E*. **77** (1): 011901. arXiv:0707.3642 (<https://arxiv.org/abs/0707.3642>). Bibcode:2008PhRvE..77a1901R (<https://ui.adsabs.harvard.edu/abs/2008PhRvE..77a1901R>). doi:10.1103/physreve.77.011901 (<https://doi.org/10.1103%2Fphysreve.77.011901>). PMID 18351870 (<https://pubmed.ncbi.nlm.nih.gov/18351870>).
- Wells, W.M. III; Viola, P.; Atsumi, H.; Nakajima, S.; Kikinis, R. (1996). "Multi-modal volume registration by maximization of mutual information" (<https://web.archive.org/web/20080906201633/http://www.ai.mit.edu/people/sw/papers/mia.pdf>) (PDF). *Medical Image Analysis*. **1** (1): 35–51. doi:10.1016/S1361-8415(01)80004-9 (<https://doi.org/10.1016%2FS1361-8415%2801%2980004-9>). PMID 9873920 (<https://pubmed.ncbi.nlm.nih.gov/9873920>). Archived from the original (<http://www.ai.mit.edu/people/sw/papers/mia.pdf>) (PDF) on 2008-09-06. Retrieved 2010-08-05.
- Pandey, Biswajit; Sarkar, Suman (2017). "How much a galaxy knows about its large-scale environment?: An information theoretic perspective". *Monthly Notices of the Royal Astronomical Society Letters*. **467** (1): L6. arXiv:1611.00283 (<https://arxiv.org/abs/1611.00283>). Bibcode:2017MNRAS.467L...6P (<https://ui.adsabs.harvard.edu/abs/2017MNRAS.467L...6P>). doi:10.1093/mnras/slw250 (<https://doi.org/10.1093%2Fmnras%2Fslw250>).

Retrieved from "https://en.wikipedia.org/w/index.php?title=Mutual_information&oldid=952431409"

This page was last edited on 22 April 2020, at 05:45 (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.