

Kullback–Leibler divergence

Not to be confused with **divergence** in calculus.

In probability theory and information theory, the **Kullback–Leibler divergence**^{[1][2][3]} (also **information divergence**, **information gain**, **relative entropy**, or **KLIC**; here abbreviated as KL divergence) is a non-symmetric measure of the difference between two probability distributions P and Q . Specifically, the Kullback–Leibler divergence of Q from P , denoted $DKL(P||Q)$, is a measure of the information lost when Q is used to approximate P :^[4] The KL divergence measures the expected number of *extra* (so intuitively, it is non negative; and can be verified by **Jensen’s inequality**) bits required to code samples from P when using a code based on Q , rather than using the true code based on P . Typically P represents the “true” distribution of data, observations, or a precisely calculated theoretical distribution. The measure Q typically represents a theory, model, description, or approximation of P .

Although it is often intuited as a **metric** or **distance**, the KL divergence is not a true **metric** — for example, it is not symmetric: the KL divergence from P to Q is generally not the same as that from Q to P . However, its infinitesimal form, specifically its **Hessian**, is a **metric tensor**: it is the **Fisher information metric**.

KL divergence is a special case of a broader class of **divergences** called f -**divergences**. It was originally introduced by **Solomon Kullback** and **Richard Leibler** in 1951 as the **directed divergence** between two distributions. It can be derived from a **Bregman divergence**.

1 Definition

For discrete probability distributions P and Q , the KL divergence of Q from P is defined to be

$$D_{KL}(P||Q) = \sum_i P(i) \ln \frac{P(i)}{Q(i)}.$$

In words, it is the expectation of the logarithmic difference between the probabilities P and Q , where the expectation is taken using the probabilities P . The KL divergence is only defined if P and Q both sum to 1 and if $Q(i) = 0$ implies $P(i) = 0$ for all i (absolute continuity). If the quantity $0 \ln 0$ appears in the formula, it is interpreted as zero because $\lim_{x \rightarrow 0} x \ln(x) = 0$.

For distributions P and Q of a continuous random variable, KL divergence is defined to be the integral:^[5]

$$D_{KL}(P||Q) = \int_{-\infty}^{\infty} p(x) \ln \frac{p(x)}{q(x)} dx,$$

where p and q denote the densities of P and Q .

More generally, if P and Q are probability measures over a set X , and P is **absolutely continuous** with respect to Q , then the Kullback–Leibler divergence from P to Q is defined as

$$D_{KL}(P||Q) = \int_X \ln \frac{dP}{dQ} dP,$$

where $\frac{dP}{dQ}$ is the **Radon–Nikodym derivative** of P with respect to Q , and provided the expression on the right-hand side exists. Equivalently, this can be written as

$$D_{KL}(P||Q) = \int_X \ln \left(\frac{dP}{dQ} \right) \frac{dP}{dQ} dQ,$$

which we recognize as the entropy of P relative to Q . Continuing in this case, if μ is any measure on X for which $p = \frac{dP}{d\mu}$ and $q = \frac{dQ}{d\mu}$ exist, then the KL divergence from P to Q is given as

$$D_{KL}(P||Q) = \int_X p \ln \frac{p}{q} d\mu.$$

The logarithms in these formulae are taken to base 2 if information is measured in units of **bits**, or to base e if information is measured in **nats**. Most formulas involving the KL divergence hold irrespective of log base.

Various conventions exist for referring to $DKL(P||Q)$ in words. Often it is referred to as the divergence *between* P and Q ; however this fails to convey the fundamental asymmetry in the relation. Sometimes it may be found described as the divergence of P from, or with respect to Q (often in the context of relative entropy, or information gain). However, in the present article the divergence of Q from P will be the language used, as this best relates to the idea that it is P that is considered the underlying “true” or “best guess” distribution, that expectations will be calculated with reference to, while Q is some divergent, less good, approximate distribution.

2 Motivation

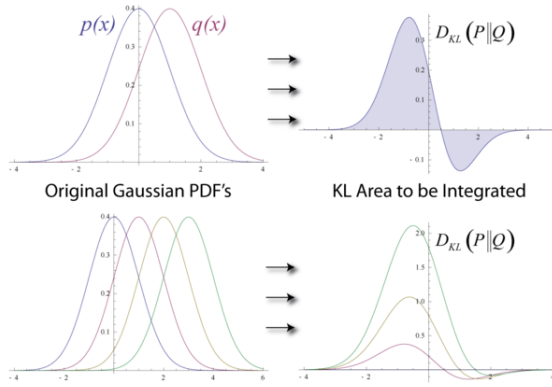


Illustration of the Kullback–Leibler (KL) divergence for two normal Gaussian distributions. Note the typical asymmetry for the KL divergence is clearly visible.

In information theory, the **Kraft–McMillan theorem** establishes that any directly decodable coding scheme for coding a message to identify one value x_i out of a set of possibilities X can be seen as representing an implicit probability distribution $q(x_i) = 2^{-l_i}$ over X , where l_i is the length of the code for x_i in bits. Therefore, KL divergence can be interpreted as the expected extra message-length per datum that must be communicated if a code that is optimal for a given (wrong) distribution Q is used, compared to using a code based on the true distribution P .

$$\begin{aligned} D_{\text{KL}}(P||Q) &= - \sum_x p(x) \log q(x) + \sum_x p(x) \log p(x) = \int_{x_a}^{x_b} P(x) \log \left(\frac{P(x)}{Q(x)} \right) dx = \\ &= \frac{H(P, Q)}{H(P, Q)} - \frac{H(P)}{H(P)} = \int_{y_a}^{y_b} P(y) \log \left(\frac{P(y) dy/dx}{Q(y) dy/dx} \right) dy = \\ &= \int_{y_a}^{y_b} P(y) \log \left(\frac{P(y)}{Q(y)} \right) dy \end{aligned}$$

where $H(P, Q)$ is called the **cross entropy** of P and Q , and $H(P)$ is the **entropy** of P .

Note also that there is a relation between the KL divergence and the "rate function" in the theory of large deviations.^{[6][7]}

Kullback brings together all notions of information in his historic text, *Information Theory and Statistics*. For instance he shows that the mean discriminating information between two hypotheses is the basis for all of the various measures of information, from Shannon to Fisher. Shannon's rate is the mean information between the hypotheses of dependence and independence of processes. Fisher's information is second order term and dominant in the Taylor approximation of the discriminating information between two models of the same parametric family.^[8]

3 Computing the closed form

For many common families of distributions, the KL divergence between two distributions in the family can be

derived in closed form. This can often be done most easily using the form of the KL divergence in terms of **expected values** or in terms of **information entropy**:

$$D_{\text{KL}}(P||Q) = -E(\ln q(x)) + E(\ln p(x)) = H(P, Q) - H(P)$$

where $H(P) = -E(\ln p(x))$ is the information entropy of P , and $H(P, Q)$ is the **cross entropy** of P and Q .

4 Properties

The Kullback–Leibler divergence is always non-negative,

$$D_{\text{KL}}(P||Q) \geq 0,$$

a result known as **Gibbs' inequality**, with $D_{\text{KL}}(P||Q)$ zero if and only if $P = Q$ **almost everywhere**. The entropy $H(P)$ thus sets a minimum value for the cross-entropy $H(P, Q)$, the expected number of bits required when using a code based on Q rather than P ; and the KL divergence therefore represents the expected number of extra bits that must be transmitted to identify a value x drawn from X , if a code is used corresponding to the probability distribution Q , rather than the "true" distribution P .

The Kullback–Leibler divergence remains well-defined for continuous distributions, and furthermore is invariant under parameter transformations. For example, if a transformation is made from variable x to variable $y(x)$, then, since $P(x) dx = P(y) dy$ and $Q(x) dx = Q(y) dy$ the Kullback–Leibler divergence may be rewritten:

where $y_a = y(x_a)$ and $y_b = y(x_b)$. Although it was assumed that the transformation was continuous, this need not be the case. This also shows that the Kullback–Leibler divergence produces a **dimensionally consistent** quantity, since if x is a dimensioned variable, $P(x)$ and $Q(x)$ are also dimensioned, since e.g. $P(x) dx$ is dimensionless. The argument of the logarithmic term is and remains dimensionless, as it must. It can therefore be seen as in some ways a more fundamental quantity than some other properties in information theory^[9] (such as **self-information** or **Shannon entropy**), which can become undefined or negative for non-discrete probabilities.

The Kullback–Leibler divergence is additive for independent distributions in much the same way as Shannon entropy. If P_1, P_2 are independent distributions, with the joint distribution $P(x, y) = P_1(x)P_2(y)$, and Q, Q_1, Q_2 likewise, then

$$D_{\text{KL}}(P||Q) = D_{\text{KL}}(P_1||Q_1) + D_{\text{KL}}(P_2||Q_2).$$

5 KL divergence for the normal distributions

The Kullback–Leibler divergence between two multivariate normal distributions of the dimension k with the means μ_0, μ_1 and their corresponding nonsingular covariance matrices Σ_0, Σ_1 is:

$$D_{\text{KL}}(\mathcal{N}_0 \parallel \mathcal{N}_1) = \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^\top \Sigma_1^{-1} (\mu_1 - \mu_0) - k - \ln \left(\frac{|\Sigma_1|}{|\Sigma_0|} \right) \right) \quad [10]$$

The logarithm in the last term must be taken to base e since all terms apart from the last are base- e logarithms of expressions that are either factors of the density function or otherwise arise naturally. The equation therefore gives a result measured in nats. Dividing the entire expression above by $\log_e 2$ yields the divergence in bits.

6 Relation to metrics

One might be tempted to call it a "distance metric" on the space of probability distributions, but this would not be correct as the Kullback–Leibler divergence is not symmetric – that is, $D_{\text{KL}}(P \parallel Q) \neq D_{\text{KL}}(Q \parallel P)$, – nor does it satisfy the triangle inequality. Still, being a premetric, it generates a topology on the space of generalized probability distributions, of which probability distributions proper are a special case. More concretely, if $\{P_1, P_2, \dots\}$ is a sequence of distributions such that

$$\lim_{n \rightarrow \infty} D_{\text{KL}}(P_n \parallel Q) = 0$$

then it is said that $P_n \xrightarrow{D} Q$. Pinsker's inequality entails that $P_n \xrightarrow{D} P \Rightarrow P_n \xrightarrow{\text{TV}} P$, where the latter stands for the usual convergence in total variation.

Following Rényi (1970, 1961)^{[11][12]} the term is sometimes also called the **information gain** about X achieved if P can be used instead of Q . It is also called the **relative entropy**, for using Q instead of P .

6.1 Fisher information metric

However, the Kullback–Leibler divergence is rather directly related to a metric, specifically, the **Fisher information metric**. This can be made explicit as follows. Assume that the probability distributions P and Q are both parameterized by some (possibly multi-dimensional) parameter θ . Consider then two close by values of $P = P(\theta)$ and $Q = P(\theta_0)$ so that the parameter θ differs by only a small amount from the parameter value θ_0 . Specifically,

up to first order one has (using the Einstein summation convention)

$$P(\theta) = P(\theta_0) + \Delta\theta^j P_j(\theta_0) + \dots$$

with $\Delta\theta^j = (\theta - \theta_0)^j$ a small change of θ in the j direction, and $P_j(\theta_0) = \frac{\partial P}{\partial \theta^j}(\theta_0)$ the corresponding rate of change in the probability distribution. Since the KL divergence has an absolute minimum 0 for $P = Q$, i.e. $\theta = \theta_0$, it changes only to *second* order in the small parameters $\Delta\theta^j$. More formally, as for any minimum, the first derivatives of the divergence vanish

$$\left. \frac{\partial}{\partial \theta^j} D_{\text{KL}}(P(\theta) \parallel P(\theta_0)) \right|_{\theta=\theta_0} = 0,$$

and by the Taylor expansion one has up to second order

$$D_{\text{KL}}(P(\theta) \parallel P(\theta_0)) = \frac{1}{2} \Delta\theta^j \Delta\theta^k g_{jk}(\theta_0) + \dots$$

where the Hessian matrix of the divergence

$$g_{jk}(\theta_0) = \left. \frac{\partial^2}{\partial \theta^j \partial \theta^k} D_{\text{KL}}(P(\theta) \parallel P(\theta_0)) \right|_{\theta=\theta_0}$$

must be positive semidefinite. Letting θ_0 vary (and dropping the subindex 0) the Hessian $g_{jk}(\theta)$ defines a (possibly degenerate) Riemannian metric on the θ parameter space, called the Fisher information metric.

7 Relation to other quantities of information theory

Many of the other quantities of information theory can be interpreted as applications of the KL divergence to specific cases.

The self-information,

$$I(m) = D_{\text{KL}}(\delta_{im} \parallel \{p_i\}),$$

is the KL divergence of the probability distribution $P(i)$ from a Kronecker delta representing certainty that $i = m$ — i.e. the number of extra bits that must be transmitted to identify i if only the probability distribution $P(i)$ is available to the receiver, not the fact that $i = m$.

The mutual information,

$$\begin{aligned} I(X; Y) &= D_{\text{KL}}(P(X, Y) \parallel P(X)P(Y)) \\ &= \mathbb{E}_X \{ D_{\text{KL}}(P(Y|X) \parallel P(Y)) \} \\ &= \mathbb{E}_Y \{ D_{\text{KL}}(P(X|Y) \parallel P(X)) \} \end{aligned}$$

is the KL divergence of the product $P(X)P(Y)$ of the two marginal probability distributions from the joint probability distribution $P(X, Y)$ — i.e. the expected number of extra bits that must be transmitted to identify X and Y if they are coded using only their marginal distributions instead of the joint distribution. Equivalently, if the joint probability $P(X, Y)$ is known, it is the expected number of extra bits that must on average be sent to identify Y if the value of X is not already known to the receiver.

The **Shannon entropy**,

$$\begin{aligned} H(X) &= (i) \mathbb{E}_x\{I(x)\} \\ &= (ii) \log N - D_{\text{KL}}(P(X) \| P_U(X)) \end{aligned}$$

is the number of bits which would have to be transmitted to identify X from N equally likely possibilities, *less* the KL divergence of the uniform distribution $P_U(X)$ from the true distribution $P(X)$ — i.e. *less* the expected number of bits saved, which would have had to be sent if the value of X were coded according to the uniform distribution $P_U(X)$ rather than the true distribution $P(X)$.

The **conditional entropy**,

$$\begin{aligned} H(X | Y) &= \log N - D_{\text{KL}}(P(X, Y) \| P_U(X)P(Y)) \\ &= (i) \log N - D_{\text{KL}}(P(X, Y) \| P(X)P(Y)) \\ &= H(X) - I(X; Y) \\ &= (ii) \log N - \mathbb{E}_Y\{D_{\text{KL}}(P(X|Y) \| P_U(X))\} \end{aligned}$$

is the number of bits which would have to be transmitted to identify X from N equally likely possibilities, *less* the KL divergence of the product distribution $P_U(X)P(Y)$ from the true joint distribution $P(X, Y)$ — i.e. *less* the expected number of bits saved which would have had to be sent if the value of X were coded according to the uniform distribution $P_U(X)$ rather than the conditional distribution $P(X | Y)$ of X given Y .

The **cross entropy** between two probability distributions measures the average number of bits needed to identify an event from a set of possibilities, if a coding scheme is used based on a given probability distribution q , rather than the “true” distribution p . The cross entropy for two distributions p and q over the same probability space is thus defined as follows:

$$H(p, q) = \mathbb{E}_p[-\log q] = H(p) + D_{\text{KL}}(p \| q).$$

8 KL divergence and Bayesian updating

In **Bayesian statistics** the KL divergence can be used as a measure of the information gain in moving from a prior distribution to a posterior distribution. If some new fact

$Y = y$ is discovered, it can be used to update the probability distribution for X from $p(x | I)$ to a new posterior probability distribution $p(x | y, I)$ using **Bayes’ theorem**:

$$p(x | y, I) = \frac{p(y | x, I)p(x | I)}{p(y | I)}$$

This distribution has a new entropy

$$H(p(\cdot | y, I)) = \sum_x p(x | y, I) \log p(x | y, I),$$

which may be less than or greater than the original entropy $H(p(\cdot | I))$. However, from the standpoint of the new probability distribution one can estimate that to have used the original code based on $p(x | I)$ instead of a new code based on $p(x | y, I)$ would have added an expected number of bits

$$D_{\text{KL}}(p(\cdot | y, I) \| p(\cdot | I)) = \sum_x p(x | y, I) \log \frac{p(x | y, I)}{p(x | I)}$$

to the message length. This therefore represents the amount of useful information, or information gain, about X , that we can estimate has been learned by discovering $Y = y$.

If a further piece of data, $Y_2 = y_2$, subsequently comes in, the probability distribution for x can be updated further, to give a new best guess $p(x | y_1, y_2, I)$. If one reinvestigates the information gain for using $p(x | y_1, I)$ rather than $p(x | I)$, it turns out that it may be either greater or less than previously estimated:

$$\begin{aligned} \sum_x p(x | y_1, y_2, I) \log \frac{p(x | y_1, y_2, I)}{p(x | I)} &\text{ may be } \leq \\ \text{or } > \text{ than } \sum_x p(x | y_1, I) \log \frac{p(x | y_1, I)}{p(x | I)} \end{aligned}$$

and so the combined information gain does *not* obey the triangle inequality:

$$\begin{aligned} D_{\text{KL}}(p(\cdot | y_1, y_2, I) \| p(\cdot | I)) &\text{ may be } <, = \\ \text{or } > \text{ than } D_{\text{KL}}(p(\cdot | y_1, y_2, I) \| p(\cdot | y_1, I)) &+ \\ D_{\text{KL}}(p(\cdot | y_1, I) \| p(\cdot | I)) \end{aligned}$$

All one can say is that on *average*, averaging using $p(y_2 | y_1, x, I)$, the two sides will average out.

8.1 Bayesian experimental design

A common goal in **Bayesian experimental design** is to maximise the expected KL divergence between the prior and the posterior.^[13] When posteriors are approximated to be Gaussian distributions, a design maximising the expected KL divergence is called **Bayes d-optimal**.

9 Discrimination information

The Kullback–Leibler divergence $DKL(p(x|H_1) \parallel p(x|H_0))$ can also be interpreted as the expected **discrimination information** for H_1 over H_0 : the mean information per sample for discriminating in favor of a hypothesis H_1 against a hypothesis H_0 , when hypothesis H_1 is true.^[14] Another name for this quantity, given to it by I.J. Good, is the expected **weight of evidence** for H_1 over H_0 to be expected from each sample.

The expected weight of evidence for H_1 over H_0 is **not** the same as the information gain expected per sample about the probability distribution $p(H)$ of the hypotheses,

$$DKL(p(x|H_1) \parallel p(x|H_0)) \neq IG = DKL(p(H|x) \parallel p(H)).$$

Either of the two quantities can be used as a **utility function** in Bayesian experimental design, to choose an optimal next question to investigate: but they will in general lead to rather different experimental strategies.

On the entropy scale of *information gain* there is very little difference between near certainty and absolute certainty—coding according to a near certainty requires hardly any more bits than coding according to an absolute certainty. On the other hand, on the **logit** scale implied by weight of evidence, the difference between the two is enormous – infinite perhaps; this might reflect the difference between being almost sure (on a probabilistic level) that, say, the **Riemann hypothesis** is correct, compared to being certain that it is correct because one has a mathematical proof. These two different scales of **loss function** for uncertainty are *both* useful, according to how well each reflects the particular circumstances of the problem in question.

9.1 Principle of minimum discrimination information

The idea of Kullback–Leibler divergence as discrimination information led Kullback to propose the Principle of **Minimum Discrimination Information** (MDI): given new facts, a new distribution f should be chosen which is as hard to discriminate from the original distribution f_0 as possible; so that the new data produces as small an information gain $DKL(f \parallel f_0)$ as possible.

For example, if one had a prior distribution $p(x,a)$ over x and a , and subsequently learnt the true distribution of a was $u(a)$, the Kullback–Leibler divergence between the new joint distribution for x and a , $q(x|a)u(a)$, and the earlier prior distribution would be:

$$DKL(q(x|a)u(a) \parallel p(x,a)) = \mathbb{E}_{u(a)}\{DKL(q(x|a) \parallel p(x|a))\} + DKL(u(a) \parallel p(a))$$

i.e. the sum of the KL divergence of $p(a)$ the prior distribution for a from the updated distribution $u(a)$, plus the

expected value (using the probability distribution $u(a)$) of the KL divergence of the prior conditional distribution $p(x|a)$ from the new conditional distribution $q(x|a)$. (Note that often the later expected value is called the *conditional KL divergence* (or *conditional relative entropy*) and denoted by $DKL(q(x|a) \parallel p(x|a))$ ^[15]) This is minimised if $q(x|a) = p(x|a)$ over the whole support of $u(a)$; and we note that this result incorporates Bayes' theorem, if the new distribution $u(a)$ is in fact a δ function representing certainty that a has one particular value.

MDI can be seen as an extension of Laplace's Principle of **Insufficient Reason**, and the Principle of **Maximum Entropy** of E.T. Jaynes. In particular, it is the natural extension of the principle of maximum entropy from discrete to continuous distributions, for which Shannon entropy ceases to be so useful (see *differential entropy*), but the KL divergence continues to be just as relevant.

In the engineering literature, MDI is sometimes called the **Principle of Minimum Cross-Entropy** (MCE) or **Minxent** for short. Minimising the KL divergence of m from p with respect to m is equivalent to minimizing the cross-entropy of p and m , since

$$H(p, m) = H(p) + DKL(p \parallel m),$$

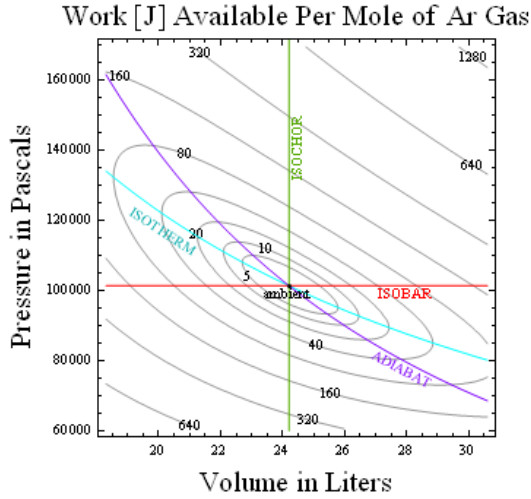
which is appropriate if one is trying to choose an adequate approximation to p . However, this is just as often *not* the task one is trying to achieve. Instead, just as often it is m that is some fixed prior reference measure, and p that one is attempting to optimise by minimising $DKL(p \parallel m)$ subject to some constraint. This has led to some ambiguity in the literature, with some authors attempting to resolve the inconsistency by redefining cross-entropy to be $DKL(p \parallel m)$, rather than $H(p, m)$.

10 Relationship to available work

Surprisals^[16] add where probabilities multiply. The surprisal for an event of probability p is defined as $s = k \ln(1/p)$. If k is $\{1, 1/\ln 2, 1.38 \times 10^{-23}\}$ then surprisal is in $\{\text{nats, bits, or } J/K\}$ so that, for instance, there are N bits of surprisal for landing all “heads” on a toss of N coins.

Best-guess states (e.g. for atoms in a gas) are inferred by maximizing the *average surprisal* S (**entropy**) for a given set of control parameters (like pressure P or volume V). This constrained **entropy maximization**, both classically^[17] and quantum mechanically,^[18] minimizes Gibbs availability in entropy units^[19] $A \equiv -k \ln Z$ where Z is a constrained multiplicity or **partition function**.

When temperature T is fixed, free energy ($T \times A$) is also minimized. Thus if T , V and number of molecules N are constant, the **Helmholtz free energy** $F \equiv U - TS$ (where U is energy) is minimized as a system “equilibrates.” If



Pressure versus volume plot of available work from a mole of Argon gas relative to ambient, calculated as T_o times KL divergence.

T and P are held constant (say during processes in your body), the **Gibbs free energy** $G = U + PV - TS$ is minimized instead. The change in free energy under these conditions is a measure of available **work** that might be done in the process. Thus available work for an ideal gas at constant temperature T_o and pressure P_o is $W = \Delta G = NkT_o\Theta(V/V_o)$ where $V_o = NkT_o/P_o$ and $\Theta(x) = x - 1 - \ln x \geq 0$ (see also **Gibbs inequality**).

More generally^[20] the **work available** relative to some ambient is obtained by multiplying ambient temperature T_o by KL divergence or **net surprisal** $\Delta I \geq 0$, defined as the average value of $k \ln(p/p_o)$ where p_o is the probability of a given state under ambient conditions. For instance, the work available in equilibrating a monatomic ideal gas to ambient values of V_o and T_o is thus $W = T_o\Delta I$, where KL divergence $\Delta I = Nk[\Theta(V/V_o) + \frac{3}{2}\Theta(T/T_o)]$. The resulting contours of constant KL divergence, shown at right for a mole of Argon at standard temperature and pressure, for example put limits on the conversion of hot to cold as in flame-powered air-conditioning or in the unpowered device to convert boiling-water to ice-water discussed here.^[21] Thus KL divergence measures thermodynamic availability in bits.

11 Quantum information theory

For density matrices P and Q on a Hilbert space the K–L divergence (or **quantum relative entropy** as it is often called in this case) from P to Q is defined to be

$$D_{\text{KL}}(P\|Q) = \text{Tr}(P(\log(P) - \log(Q))).$$

In **quantum information science** the minimum of $D_{\text{KL}}(P\|Q)$ over all separable states Q can also be used

as a measure of **entanglement** in the state P .

12 Relationship between models and reality

Just as KL divergence of “ambient from actual” measures thermodynamic availability, KL divergence of “model from reality” is also useful even if the only clues we have about reality are some experimental measurements. In the former case KL divergence describes *distance to equilibrium* or (when multiplied by ambient temperature) the amount of *available work*, while in the latter case it tells you about surprises that reality has up its sleeve or, in other words, *how much the model has yet to learn*.

Although this tool for evaluating models against systems that are accessible experimentally may be applied in any field, its application to models in ecology via **Akaike information criterion** are particularly well described in papers^[22] and a book^[23] by Burnham and Anderson. In a nutshell the KL divergence of a model from reality may be estimated, to within a constant additive term, by a function (like the squares summed) of the deviations observed between data and the model’s predictions. Estimates of such divergence for models that share the same additive term can in turn be used to choose between models.

When trying to fit parametrized models to data there are various estimators which attempt to minimize Kullback–Leibler divergence, such as **maximum likelihood** and **maximum spacing** estimators.

13 Symmetrised divergence

Kullback and Leibler themselves actually defined the divergence as:

$$D_{\text{KL}}(P\|Q) + D_{\text{KL}}(Q\|P)$$

which is symmetric and nonnegative. This quantity has sometimes been used for feature selection in **classification** problems, where P and Q are the conditional pdfs of a feature under two different classes.

An alternative is given via the λ divergence,

$$D_{\lambda}(P\|Q) = \lambda D_{\text{KL}}(P\|\lambda P + (1-\lambda)Q) + (1-\lambda) D_{\text{KL}}(Q\|\lambda P + (1-\lambda)Q),$$

which can be interpreted as the expected information gain about X from discovering which probability distribution X is drawn from, P or Q , if they currently have probabilities λ and $(1-\lambda)$ respectively.

The value $\lambda = 0.5$ gives the **Jensen–Shannon divergence**, defined by

$$D_{JS} = \frac{1}{2}D_{KL}(P||M) + \frac{1}{2}D_{KL}(Q||M)$$

where M is the average of the two distributions,

$$M = \frac{1}{2}(P + Q).$$

DJS can also be interpreted as the capacity of a noisy information channel with two inputs giving the output distributions p and q . The Jensen–Shannon divergence, like all f-divergences, is *locally* proportional to the **Fisher information metric**. It is similar to the Hellinger metric (in the sense that induces the same affine connection on a statistical manifold), and equal to one-half the so-called *Jeffreys divergence* (Rubner et al., 2000; Jeffreys 1946^[24]).

14 Relationship to Rényi divergence and Hellinger distance

If P and Q are two probability measures, then the squared **Hellinger distance** is the quantity given by

$$H^2(P, Q) = \frac{1}{2} \int \left(\sqrt{\frac{dP}{d\lambda}} - \sqrt{\frac{dQ}{d\lambda}} \right)^2 d\lambda$$

and the **Rényi divergence** of order α is

$$D_\alpha(P||Q) = \frac{1}{\alpha - 1} \log \int \left(\frac{dP}{d\lambda} \right)^\alpha \left(\frac{dQ}{d\lambda} \right)^{1-\alpha} d\lambda$$

The Kullback–Leibler divergence can be lower bounded in terms of the **Rényi divergence** of order $1/2$ ^[25] and the **Hellinger distance**^[26]

$$D_{KL}(P||Q) \geq D_{1/2}(P||Q) \geq 2H^2(P, Q).$$

15 Other probability-distance measures

Other measures of probability distance are the *histogram intersection*, *Chi-squared statistic*, *quadratic form distance*, *match distance*, *Kolmogorov–Smirnov distance*, and *earth mover’s distance*.^[27]

16 Data differencing

Main article: **Data differencing**

Just as *absolute* entropy serves as theoretical background for **data compression**, *relative* entropy serves as theoretical background for **data differencing** – the absolute entropy of a set of data in this sense being the data required to reconstruct it (minimum compressed size), while the relative entropy of a target set of data, given a source set of data, is the data required to reconstruct the target *given* the source (minimum size of a patch).

17 See also

- **Bregman divergence**
- **Jensen–Shannon divergence**
- **Deviance information criterion**
- **Bayesian information criterion**
- **Quantum relative entropy**
- **Information gain in decision trees**
- **Solomon Kullback and Richard Leibler**
- **Information theory and measure theory**
- **Entropy power inequality**
- **Information gain ratio**
- **Entropic value at risk**
- **Akaike Information Criterion**

18 References

- [1] Kullback, S.; Leibler, R.A. (1951). “On Information and Sufficiency”. *Annals of Mathematical Statistics* **22** (1): 79–86. doi:10.1214/aoms/1177729694. MR 39968.
- [2] S. Kullback (1959) *Information theory and statistics* (John Wiley and Sons, NY).
- [3] Kullback, S. (1987). “Letter to the Editor: The Kullback–Leibler distance”. *The American Statistician* **41** (4): 340–341. JSTOR 2684769.
- [4] Kenneth P. Burnham, David R. Anderson (2002), *Model Selection and Multi-Model Inference: A Practical Information-Theoretic Approach*. Springer. (2nd ed), p.51
- [5] C. Bishop (2006). *Pattern Recognition and Machine Learning*. p. 55.
- [6] Sanov I.N. (1957) “On the probability of large deviations of random magnitudes”. *Matem. Sbornik*, v. 42 (84), 11–44.
- [7] Novak S.Y. (2011) ch. 14.5, “Extreme value methods with applications to finance”. Chapman & Hall/CRC Press. ISBN 978-1-4398-3574-6.

- [8] Kullback(1959), Information Theory and Statistics, Dover Press. ISBN 0-486-69684-7.
- [9] See the section “differential entropy - 4” in [Relative Entropy](#) video lecture by Sergio Verdú NIPS 2009
- [10] J. Duchi, Derivations for Linear Algebra and Optimization . pp. 13
- [11] A. Rényi (1970). *Probability Theory*. New York: Elsevier. Appendix, Sec.4. ISBN 0-486-45867-9.
- [12] A. Rényi (1961). “On measures of information and entropy”. “Proceedings of the 4th Berkeley Symposium on Mathematics, Statistics and Probability 1960”. pp. 547–561.
- [13] Chaloner K. and Verdinelli I. (1995) Bayesian Experimental Design: A Review. *Statistical Science* **10** (3): 273–304. doi:10.1214/ss/1177009939
- [14] Press, WH; Teukolsky, SA; Vetterling, WT; Flannery, BP (2007). “Section 14.7.2. Kullback–Leibler Distance”. *Numerical Recipes: The Art of Scientific Computing* (3rd ed.). New York: Cambridge University Press. ISBN 978-0-521-88068-8
- [15] Thomas M. Cover, Joy A. Thomas (1991) *Elements of Information Theory* (John Wiley and Sons, New York, NY), p.22
- [16] Myron Tribus (1961) *Thermodynamics and thermostatics* (D. Van Nostrand, New York)
- [17] E. T. Jaynes (1957) Information theory and statistical mechanics, *Physical Review* **106**:620
- [18] E. T. Jaynes (1957) Information theory and statistical mechanics II, *Physical Review* **108**:171
- [19] J.W. Gibbs (1873) A method of geometrical representation of thermodynamic properties of substances by means of surfaces, reprinted in *The Collected Works of J. W. Gibbs, Volume I Thermodynamics*, ed. W. R. Longley and R. G. Van Name (New York: Longmans, Green, 1931) footnote page 52.
- [20] M. Tribus and E. C. McIrvine (1971) Energy and information, *Scientific American* **224**:179–186.
- [21] P. Fraundorf (2007) Thermal roots of correlation-based complexity, *Complexity* **13**:3, 18–26
- [22] Kenneth P. Burnham and David R. Anderson (2001) Kullback–Leibler information as a basis for strong inference in ecological studies, *Wildlife Research* **28**:111–119.
- [23] Burnham, K. P. and Anderson D. R. (2002) *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach, Second Edition* (Springer Science, New York) ISBN 978-0-387-95364-9.
- [24] H. Jeffreys (1945) An Invariant Form for the Prior Probability in Estimation Problems
- [25] Van Erven, Tim; Harremoës, Peter (2014). “Rényi Divergence and Kullback-Leibler Divergence”. *IEEE Transactions on Information Theory* **60** (7): 3797–3820. doi:10.1109/TIT.2014.2320500.
- [26]
- [27] Rubner, Y., Tomasi, C., and Guibas, L. J., 2000. The Earth Mover’s distance as a metric for image retrieval. *International Journal of Computer Vision*, **40**(2): 99–121.

19 External links

- [Information Theoretical Estimators Toolbox](#)
- [Ruby gem for calculating KL divergence](#)
- [Jon Shlens’ tutorial on Kullback–Leibler divergence and likelihood theory](#)
- [Matlab code for calculating KL divergence for discrete distributions](#)
- [Sergio Verdú, Relative Entropy, NIPS 2009. One-hour video lecture.](#)
- [A modern summary of info-theoretic divergence measures](#)

20 Text and image sources, contributors, and licenses

20.1 Text

- **Kullback–Leibler divergence** *Source:* http://en.wikipedia.org/wiki/Kullback–Leibler_divergence?oldid=631351032 *Contributors:* Edward, Michael Hardy, Den fjättrade ankan, Kevin Baas, Cyan, BenKovitz, Charles Matthews, Jitse Niesen, Mpost89, Benwing, Mottzo, Wile E. Heresiarch, Giftlite, Inkling, Romanpoet, Dmb000006, Stern, Schizoid, MarkSweep, AndrewKeenanRichardson, MisterSheik, 3mta3, PAR, Jheald, Oleg Alexandrov, MartinSpacek, Linas, Shreevatsa, BlaiseFEgan, Rjwilmsi, Mathbot, Jmorgan, Spacepotato, Nothing1212, Mike Lin, Gzabers, Avraham, MDReid, Mebden, SmackBot, Mmernex, Adfernandes, Eskimbot, Object01, MclD, Ignacioerrico, Nbarth, Tsca.bot, Memming, Jon Awbrey, Dnavarro, Cronholm144, Nijdam, Dfass, Kyellan, Yoderj, A. Pichler, JForget, Thermochap, Physic sox, Sir Vicious, Winterfors, Neonleob, Mct mht, FilipeS, Amit Moscovich, Rkrish67, Stangaa, RogierBrussee, STBot, Wullj, Andre.holzner, Smite-Meister, LordAnubisBOT, Epistemenical, Punkstar89, TXiKiBoT, Miranda, Mundhenk, Jamelan, Loniousmonk, Forwardmeasure, Brech, Melcombe, Rinconsaleao, Wittnate, Sun Creator, Kaba3, Edg2103, Qwfp, DumZiBoT, Addbot, Deepmath, Fyrael, Wikomidia, Baisemain, Lightbot, Luckas-bot, Yobot, Legendre17, AnomieBOT, ￼, Obersachsebot, Xqbot, GrouchoBot, Nathanielvirgo, Chjoaygame, X7q, Citation bot 1, Kiefer.Wolfowitz, Stpasha, Angelorf, Amkilpatrick, RjwilmsiBot, Ereiniona, Kastchei, Cstahlhut, Slawekb, Quondum, ClueBot NG, Helpful Pixie Bot, Epomqo, SciCompTeacher, Francis liberty, M.daryalal, Iturrate, SFK2, Szzoli, Sunil.log, Zoltan szabo, Engheta, Velvel2 and Anonymous: 114

20.2 Images

- **File:ArgonKLdivergence.png** *Source:* <http://upload.wikimedia.org/wikipedia/commons/c/c2/ArgonKLdivergence.png> *License:* CC-BY-SA-3.0-2.5-2.0-1.0 *Contributors:* Own work *Original artist:* P. Fraundorf
- **File:KL-Gauss-Example.png** *Source:* <http://upload.wikimedia.org/wikipedia/en/a/a8/KL-Gauss-Example.png> *License:* CC-BY-SA-3.0 *Contributors:* T. Nathan Mundhenk, Ph.D thesis appendix C. *Original artist:* Mundhenk (talk)

20.3 Content license

- Creative Commons Attribution-Share Alike 3.0