



Entropy



Where does the term entropy come from?

- ❑ The term "entropy" was initially discovered by German physicist and mathematician Rudolf Clausius and was used in the field of thermodynamics.
 - ❑ Clausius initially defined the entropy of a thermodynamic system as the ratio of the heat absorbed by the system to its absolute temperature when the system is in thermal equilibrium.
- ❑ The concept of entropy was later expanded to information theory by Claude Shannon, an American mathematician and electrical engineer, in 1948. Shannon introduced the term "entropy" to quantify the amount of uncertainty or randomness in a message or signal.
 - ❑ Shannon was also known as the 'father of information theory' as he had invented the field of information theory with the paper of "*A Mathematical Theory of Communication*".

A Mathematical Theory of Communication By Father of Information Theory

- ❑ Claude Shannon demonstrated that all communication systems fundamentally share the same structure for transmitting and receiving information. He showed how messages can be encoded in bits to compress them and transmit them with near-perfect accuracy.

34

The Mathematical Theory of Communication

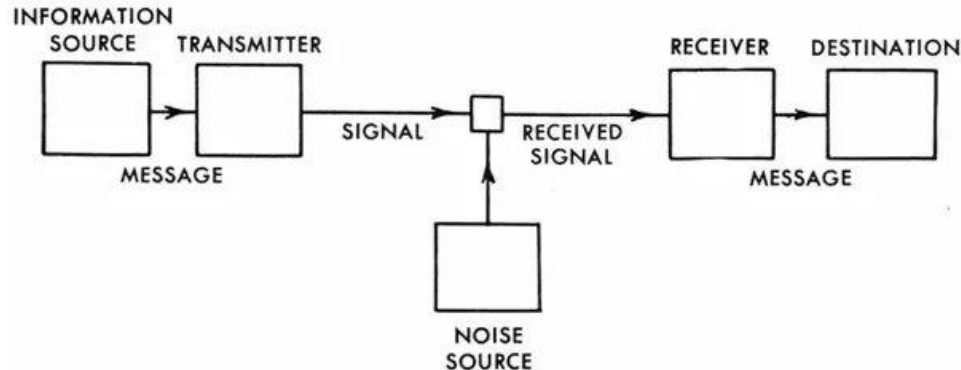
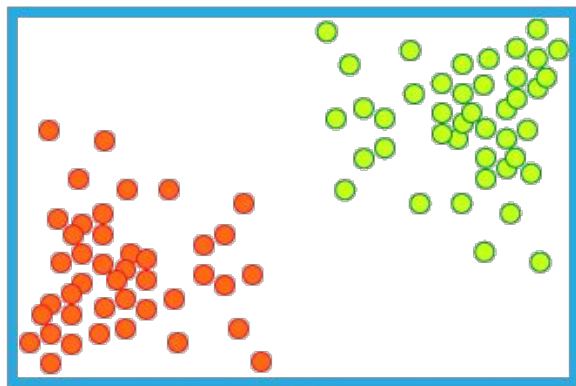


Fig. 1. — Schematic diagram of a general communication system.

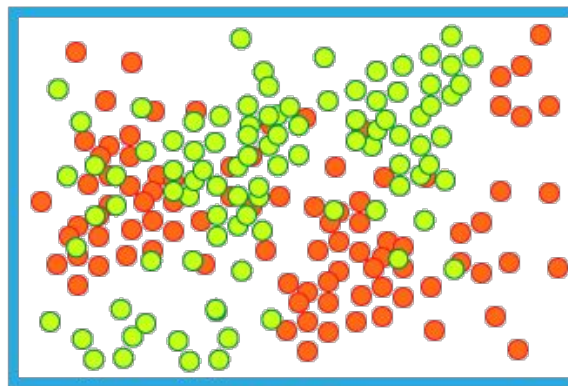
- ❑ In his article, he embarked on measuring the statistical nature of "lost information" in telephone line signals. The study aimed to address the problem of how to encode the desired information to be transmitted by a sender most effectively. For this purpose, the concept of information entropy was developed as a means to estimate the information content in a message, which is a measure of the reduced uncertainty caused by the message.

What is Information ?

- ❑ In Shannon's terms, **information entropy** is the degree of uncertainty, ie, a mathematical concept that **measures the amount of uncertainty** or randomness in a set of possible outcomes.
- ❑ In Shannon's terms, information measures surprise, or the amount of uncertainty we overcome. As a basic illustration of this point, he asked us to think about tossing a coin. A fair coin carries equal odds of landing on either side. As Shannon put it, such a coin stores one bit of information. A coin with heads on both sides carries no information.
 - ❑ An event with probability 100% is perfectly unsurprising.



Low Entropy



High Entropy

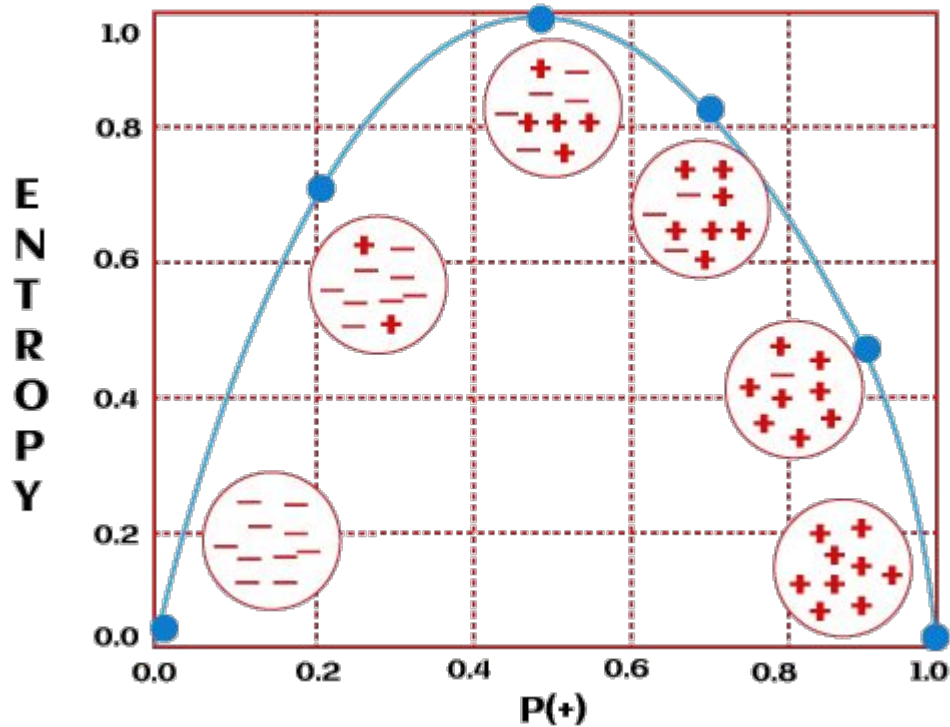
Shannon Entropy

- The entropy formula for an event X with n possible outcomes and probabilities p_1, \dots, p_n :

$$H(X) = H(p_1, \dots, p_n) = - \sum_{i=1}^n p_i \log_2 p_i$$

$$\begin{aligned} H(X) &= - \sum_{i=1}^n p(x_i) \log_2 p(x_i) \\ &= - \sum_{i=1}^2 \frac{1}{2} \log_2 \frac{1}{2} \\ &= - \sum_{i=1}^2 \frac{1}{2} \cdot (-1) = 1 \end{aligned}$$

Shannon Entropy

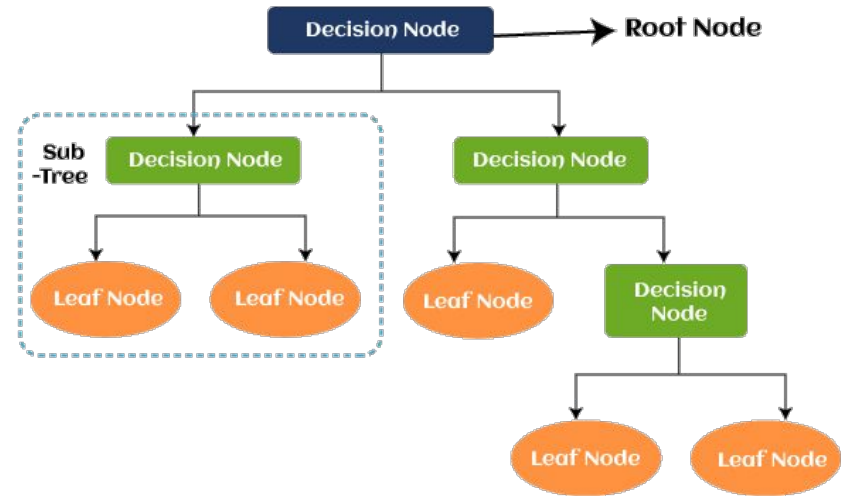


- When entropy becomes 0, then the dataset has no impurity. Datasets with 0 impurities are not useful for learning. Further, if the entropy is 1, then this kind of dataset is good for learning.

entropy is the machine learning metric that measures the unpredictability

Use Cases of Entropy in ML

- Decision trees are used to predict an outcome based on historical data. The decision tree works on the sequence of 'if-then-else' statements and a root which is our initial problem to solve.
- Entropy is used in decision trees as a measure of impurity or **disorder in a set of data**. The goal of a decision tree is to divide a dataset into subsets that are as pure (or homogeneous) as possible, and entropy is a useful tool for determining the **degree of impurity** in a set of data. When **building a decision tree**, the **algorithm will use entropy to calculate the potential information gain from each split, and will choose the split that results in the greatest reduction in entropy**. This process is repeated recursively until the tree is built.



Example of How Entropy is Used in Decision Tree Construction

1) The first step is to calculate the entropy of the target variable "Playing Golf"

So, the entropy of the dataset is: $E(\text{PlayGolf}) = -(9/14)\log_2(9/14) - (5/14)\log_2(5/14) = 0.94$

Play Golf(14)	
Yes	No
9	5

2) For the other attributes, we need to calculate the entropy after each of the split. $E(\text{PlayGolf}, \text{Outlook})$, $E(\text{PlayGolf}, \text{Temperature})$, $E(\text{PlayGolf}, \text{Humidity})$, $E(\text{PlayGolf}, \text{Windy})$

$$E(\text{Sunny}) = -(3/5)\log_2(3/5) - (2/5)\log_2(2/5) = 0.971$$

$$E(\text{Overcast}) = -(4/4)\log_2(4/4) - (0/4)\log_2(0/4) = 0$$

$$E(\text{Rainy}) = -(2/5)\log_2(2/5) - (3/5)\log_2(3/5) = 0.971$$

		PlayGolf(14)		
		Yes	No	
Outlook	Sunny	3	2	5
	Overcast	4	0	4
	Rainy	2	3	5

$$E(\text{PlayGolf}, \text{Outlook}) = P(\text{Sunny})E(\text{Sunny}) + P(\text{Overcast})E(\text{Overcast}) + P(\text{Rainy})E(\text{Rainy})$$

$$\begin{aligned}
 E(\text{PlayGolf}, \text{Outlook}) &= \frac{5}{14}E(3,2) + \frac{4}{14}E(4,0) + \frac{5}{14}E(2,3) \\
 &= \frac{5}{14}0.971 + \frac{4}{14}0.0 + \frac{5}{14}0.971 \\
 &= 0.357 * 0.971 + 0.0 + 0.357 * 0.971 \\
 &= 0.693
 \end{aligned}$$

Example of How Entropy is Used in Decision Tree Construction

3) 3rd step is to calculate the entropy for other variable called Temperature

$$E(\text{PlayGolf}, \text{Temperature}) = P(\text{Hot}) E(2,2) + P(\text{Cold}) E(3,1) + P(\text{Mild}) E(4,2)$$

$$E(\text{PlayGolf}, \text{Temperature}) = 4/14 * E(\text{Hot}) + 4/14 * E(\text{Cold}) + 6/14 * E(\text{Mild})$$

$$E(\text{PlayGolf}, \text{Temperature}) = 4/14 * E(2, 2) + 4/14 * E(3, 1) + 6/14 * E(4, 2)$$

$$\begin{aligned} E(\text{PlayGolf}, \text{Temperature}) &= 4/14 * -(2/4 \log 2/4) - (2/4 \log 2/4) \\ &+ 4/14 * -(3/4 \log 3/4) - (1/4 \log 1/4) \\ &+ 6/14 * -(4/6 \log 4/6) - (2/6 \log 2/6) \end{aligned}$$

$$\begin{aligned} E(\text{PlayGolf}, \text{Temperature}) &= 5/14 * 1.0 \\ &+ 4/14 * 1.811 \\ &+ 5/14 * 0.918 \\ &= 0.911 \end{aligned}$$

		PlayGolf(14)		
		Yes	No	
Temperature	Hot	2	2	4
	Cold	3	1	4
	Mild	4	2	6

Example of How Entropy is Used in Decision Tree Construction

4) Lets say we have the entropies for all the four attributes, let's go ahead to summarize them as shown in below:

$$E(\text{PlayGolf}, \text{Outlook}) = 0.693$$

$$E(\text{PlayGolf}, \text{Temperature}) = 0.911$$

$$E(\text{PlayGolf}, \text{Humidity}) = 0.788$$

$$E(\text{PlayGolf}, \text{Windy}) = 0.892$$

5) The next step is to calculate the information gain for each of the attributes. The attribute with the **largest information gain is used for the split**. $\text{Gain}(S, T) = \text{Entropy}(S) - \text{Entropy}(S, T)$

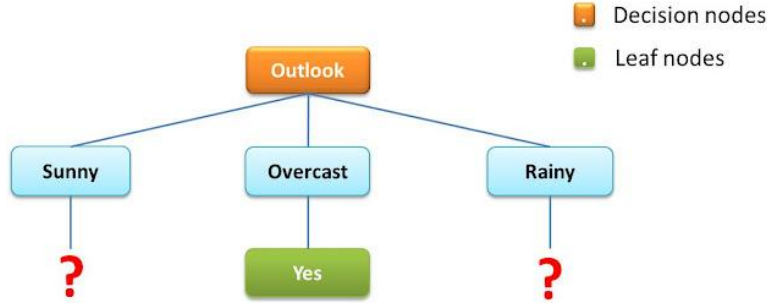
$$\text{Gain}(\text{PlayGolf}, \text{Outlook}) = \text{Entropy}(\text{PlayGolf}) - \text{Entropy}(\text{PlayGolf}, \text{Outlook}) = 0.94 - 0.693 = 0.247$$

$$\text{Gain}(\text{PlayGolf}, \text{Temperature}) = \text{Entropy}(\text{PlayGolf}) - \text{Entropy}(\text{PlayGolf}, \text{Temperature}) = 0.94 - 0.911 = 0.029$$

$$\text{Gain}(\text{PlayGolf}, \text{Humidity}) = \text{Entropy}(\text{PlayGolf}) - \text{Entropy}(\text{PlayGolf}, \text{Humidity}) = 0.94 - 0.788 = 0.152$$

$$\text{Gain}(\text{PlayGolf}, \text{Windy}) = \text{Entropy}(\text{PlayGolf}) - \text{Entropy}(\text{PlayGolf}, \text{Windy}) = 0.94 - 0.892 = 0.048$$

Example of How Entropy is Used in Decision Tree Construction



Decision Tree after first split

Outlook	Temperature	Humidity	Windy	Play Golf
Rainy	Hot	High	FALSE	No
Rainy	Hot	High	TRUE	No
Rainy	Mild	High	FALSE	No

Rainy	Cool	Normal	FALSE	Yes
Rainy	Mild	Normal	TRUE	Yes

The Rainy attribute could be split using High and Normal attributes and that would give us the tree below.

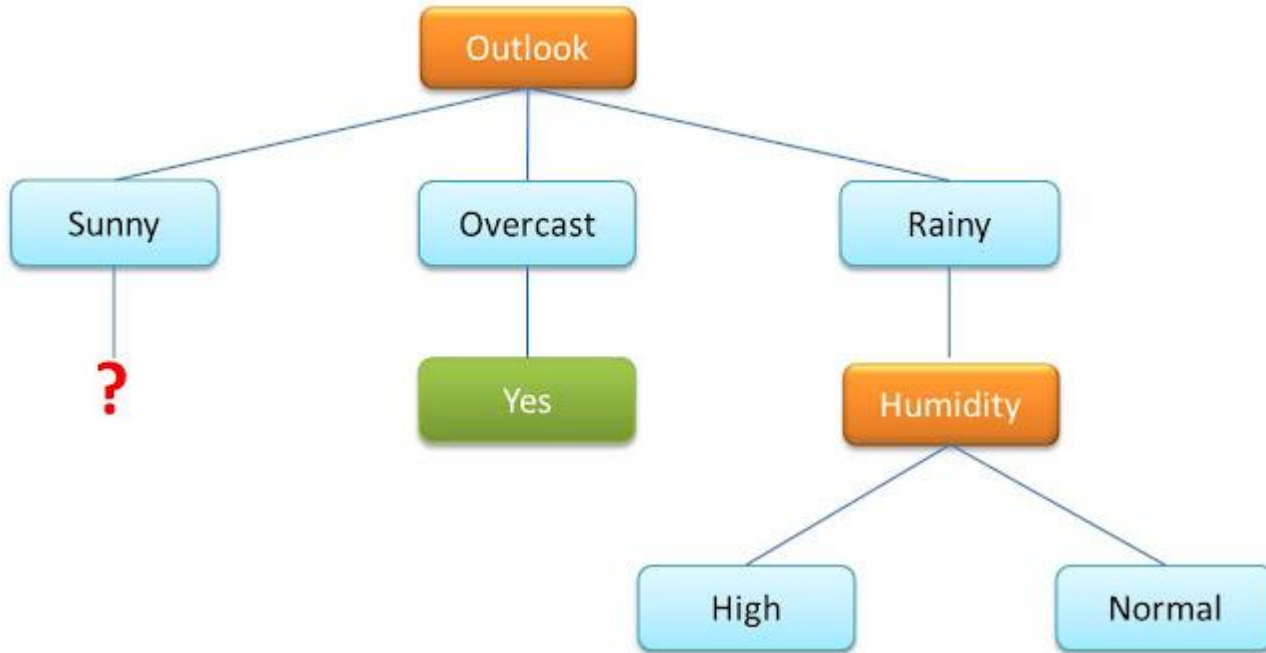
Outlook	Temperature	Humidity	Windy	Play Golf
Sunny	Mild	Normal	FALSE	Yes
Sunny	Mild	High	FALSE	Yes
Sunny	Cool	Normal	FALSE	Yes
Sunny	Cool	Normal	TRUE	No
Sunny	Mild	High	TRUE	No

Overcast	Hot	High	FALSE	Yes
Overcast	Mild	High	TRUE	Yes
Overcast	Hot	Normal	FALSE	Yes
Overcast	Cool	Normal	TRUE	Yes

Rainy	Hot	High	FALSE	No
Rainy	Hot	High	TRUE	No
Rainy	Mild	High	FALSE	No
Rainy	Cool	Normal	FALSE	Yes
Rainy	Mild	Normal	TRUE	Yes

Initial Split using Outlook : The Sunny and the Rainy attributes need to be split. The Rainy outlook can be split using either Temperature, Humidity or Windy.

Example of How Entropy is Used in Decision Tree Construction



Cross Entropy Loss (Log Loss, Logistic Loss)

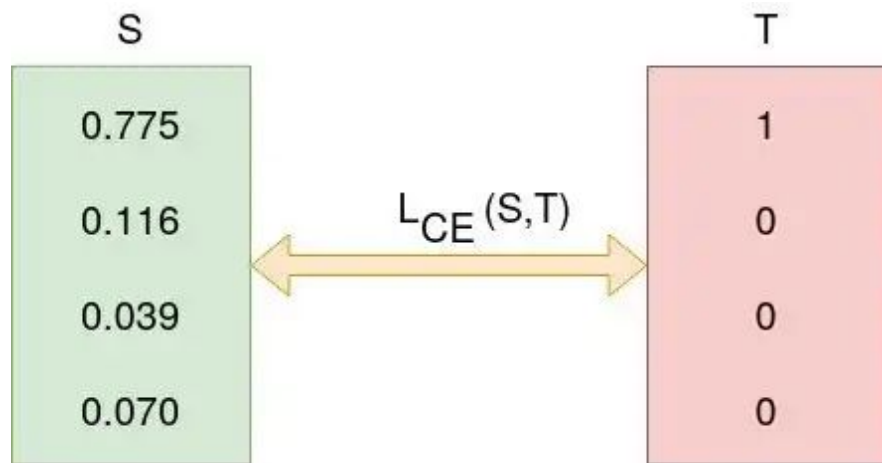
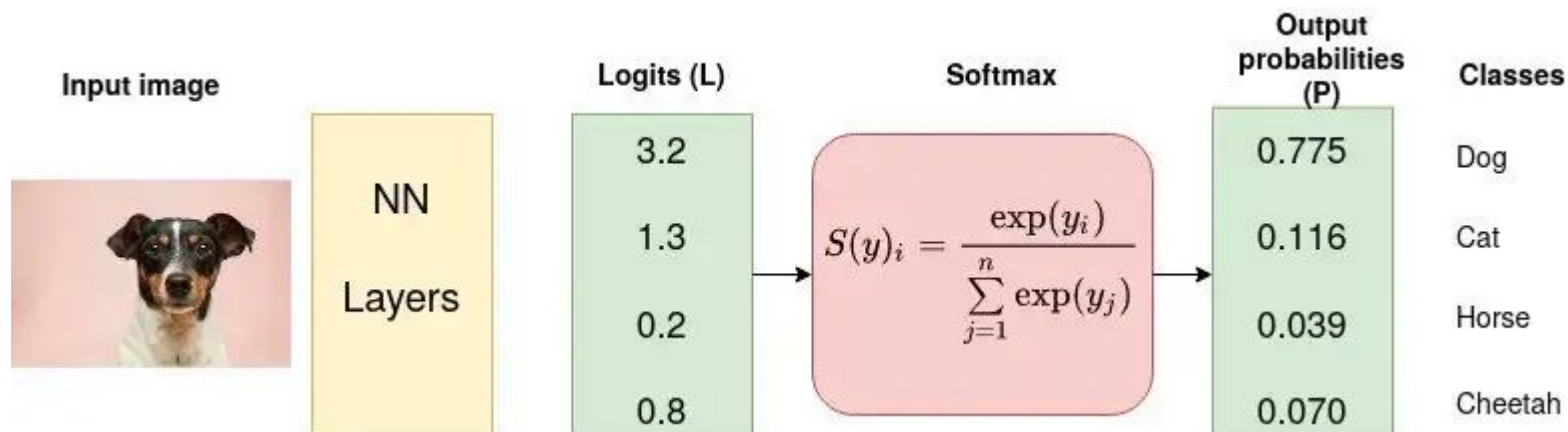
- ❑ Entropy is the number of bits required to transmit a randomly selected event from a probability distribution.
- ❑ **The cross entropy is the average number of bits needed to encode data coming from a source with distribution p when we use model q .**
- ❑ Cross-entropy is used to measure the **dissimilarity** between the predicted probability distribution and the true probability distribution. In simple terms, cross-entropy loss is a measure of the difference between the predicted class probabilities and the true class probabilities. The goal of the training process is to minimize the cross-entropy loss, which means to make the predicted class probabilities as close as possible to the true class probabilities.
- ❑ **$H(P, Q)$: Where $H()$ is the cross-entropy function, P may be the target distribution and Q is the approximation of the target distribution.** Cross-entropy can be calculated using the probabilities of the events from P and Q , as follows: $H(P, Q) = - \sum_{x \in X} P(x) * \log(Q(x))$

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

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

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

$$KL(P||Q) = \sum p_i(x) \log\left(\frac{p_i(x)}{q_i(x)}\right)$$



CROSS-ENTROPY

$S(Y)$

0.7
0.2
0.1

$D(S, L) = - \sum_i L_i \log(S_i)$

L

1.0
0.0
0.0

A still from the movie Toy Story showing Woody and Buzz Lightyear. Woody is on the left, looking slightly concerned. Buzz is on the right, looking excited and pointing his finger. The background is a simple room with a door and some toys on the floor.

INFORMATION

**INFORMATION
EVERYWHERE**

REFERENCES

[%20log\(Q\(x\)\)](https://medium.com/unpackai/cross-entropy-loss-in-ml-d9f22fc11fe0#:~:text=Cross%2Dentropy%20can%20be%20calculated,)

<https://towardsdatascience.com/entropy-is-a-measure-of-uncertainty-e2c000301c2c>

<https://www.quora.com/How-would-you-explain-Shannons-information-theory-in-laymans-terms>

https://en.wikipedia.org/wiki/Cross_entropy

https://en.wikipedia.org/wiki/Kullback%E2%80%93Leibler_divergence

<https://kindsonthegenius.com/blog/how-to-build-a-decision-tree-for-classification-step-by-step-procedure-using-entropy-and-gain/>