

Information Measures for Continuous Distributions

I-Hsiang Wang

Department of Electrical Engineering
National Taiwan University

ihwang@ntu.edu.tw

November 2, 2023

Entropy of a continuous random variable

Question: What is the entropy of a continuous real-valued r.v. X ?

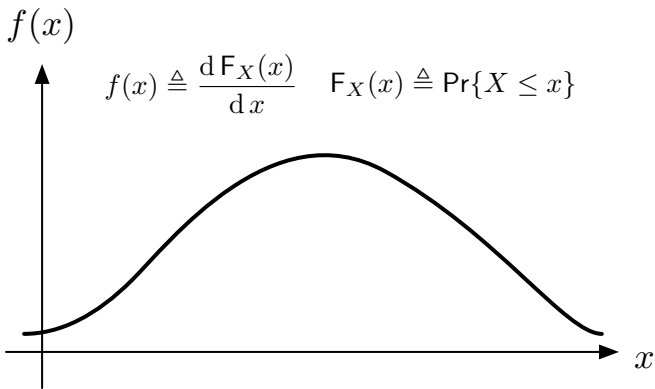
Suppose X has the probability density function (PDF) $f_X(\cdot)$.

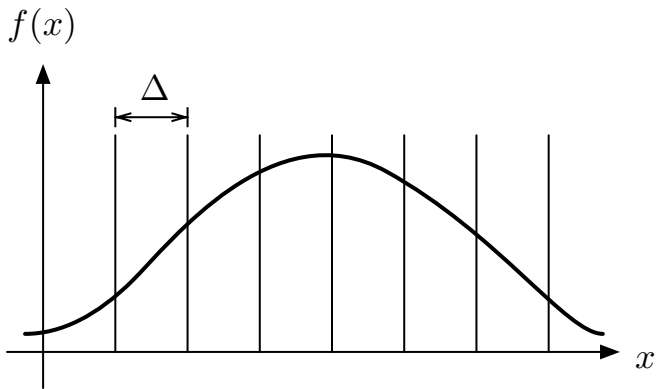
Let us discretize X to answer this question, as follows:

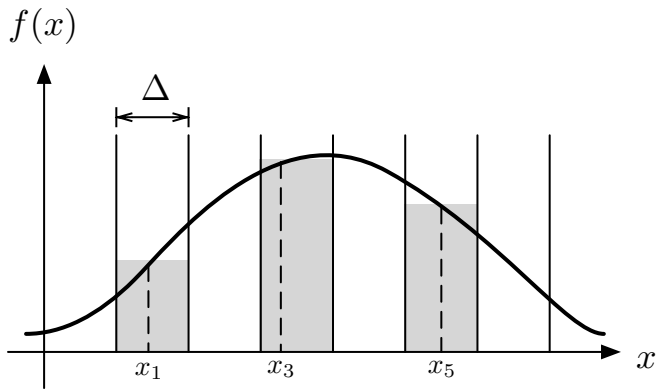
- Partition \mathbb{R} into length- Δ intervals: $\mathbb{R} = \bigcup_{k=-\infty}^{\infty} [k\Delta, (k+1)\Delta)$.
- Suppose that $f_X(\cdot) \equiv f(\cdot)$ is continuous, then by the mean-value theorem (MVT),

$$\forall k \in \mathbb{Z}, \exists x_k \in [k\Delta, (k+1)\Delta) \text{ such that } f(x_k) = \frac{1}{\Delta} \int_{k\Delta}^{(k+1)\Delta} f(x) dx.$$

- Set $[X]_{\Delta} := x_k$ if $X \in [k\Delta, (k+1)\Delta)$, with PMF $P(x_k) = f(x_k)\Delta$.







Observation: $\lim_{\Delta \rightarrow 0} H([X]_{\Delta}) = H(X)$ (intuitively), while

$$\begin{aligned} H([X]_{\Delta}) &= - \sum_{k=-\infty}^{\infty} (f(x_k)\Delta) \log(f(x_k)\Delta) \\ &= -\Delta \sum_{k=-\infty}^{\infty} f(x_k) \log f(x_k) - \log \Delta \\ &\rightarrow - \int_{-\infty}^{\infty} f(x) \log f(x) dx + \infty = \infty \quad \text{as } \Delta \rightarrow 0 \end{aligned}$$

Hence, $H(X) = \infty$ if $-\int_{-\infty}^{\infty} f(x) \log f(x) dx = \mathbb{E} \left[\log \frac{1}{f(X)} \right]$ exists.

It is quite intuitive that the entropy of a continuous random variable can be arbitrarily large, because it can take infinitely many possible values.

The term $\mathbb{E}_{X \sim f_X} [\log 1/f_X(X)]$ shares a similar form as the Shannon entropy $\mathbb{E}_{X \sim p_X} [\log 1/p_X(X)]$, and it is called the *differential entropy* of the density f_X (or the continuous random variable X).

Differential entropy

Definition 1 (Differential Entropy)

The differential entropy of a continuous r.v. X with PDF f_X is defined as

$$h(X) := \mathbb{E}_{X \sim f_X} \left[\log \frac{1}{f_X(X)} \right]$$

if the (improper) integral exists.

The previous quantization argument implies that the Shannon entropy of the quantized version of the continuous r.v. X to n -bit precision ($\Delta = 2^{-n}$) is roughly $n + h(X)$ bits.

In words, $h(X)$ is the *extra* number of bits on the average required to describe X to n -bit precision. It can be positive, zero, or negative.

Example 1 (Differential entropy of a uniform r.v.)

For a r.v. $X \sim \text{Unif}([a, b])$, that is, its PDF $f_X(x) = \frac{1}{b-a} \mathbb{1}\{a \leq x \leq b\}$, its differential entropy

$$h(X) = \log(b - a).$$

Interpretation via quantization:

- Suppose $X \sim \text{Unif}([0, 1])$. $h(X) = 0$. To describe X to n -bit precision, n bits are needed.
- Suppose $X \sim \text{Unif}([0, 2^{-m}])$, $m > 0$. $h(X) = -m$. To describe X to n -bit precision, just $n - m$ bits are needed.
- Suppose $X \sim \text{Unif}([0, 2^m])$, $m > 0$. $h(X) = m$. To describe X to n -bit precision, an extra m bits are needed.

Remark: if $b - a$ approaches 0, X becomes deterministic, and its differential entropy becomes $-\infty$.

Example 2 (Differential entropy of $N(0, 1)$)

For a r.v. $X \sim N(0, 1)$, that is, its PDF $f_X(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{x^2}{2})$, its differential entropy

$$h(X) = \frac{1}{2} \log(2\pi e).$$

Conditional differential entropy

Suppose X is jointly distributed with another random variable Y , and given $Y = y$, X is still a continuous r.v. with density $f_{X|Y}(\cdot|y)$. The differential entropy of this density is

$$h(X|Y = y) = \int_{-\infty}^{\infty} f_{X|Y}(x|y) \log \frac{1}{f_{X|Y}(x|y)} dx.$$

Averaging it over Y , the **conditional differential entropy** emerges:

$$\begin{aligned} h(X|Y) &= \begin{cases} \sum_{y \in \mathcal{Y}} p_Y(y) h(X|Y = y), & \text{if } Y \text{ is discrete with PMF } p_Y. \\ \int_{-\infty}^{\infty} f_Y(y) h(X|Y = y) dy, & \text{if } Y \text{ is continuous with PDF } f_Y. \end{cases} \\ &= E_{X,Y} \left[\log \frac{1}{f_{X|Y}(X|Y)} \right]. \end{aligned}$$

Information divergence

Recall that the information divergence of a distribution with density $f_1(\cdot)$ from another distribution with density $f_2(\cdot)$ is just

$$D(f_1 \| f_2) := E_{X \sim f_1} \left[\log \frac{f_1(X)}{f_2(X)} \right] = \int_{x \in \text{supp}_{f_1}} f_1(x) \log \frac{f_1(x)}{f_2(x)} dx$$

if the (improper) integral exists. Other properties remain to hold such as non-negativity, DPI, chain rule, conditioning increases divergence, etc..

Differential entropy turns out to have a strong connection to KL divergence. In $D(f \| g)$, suppose we replace the second PDF g by the “uniform” density $\mathbb{1}\{x \in \mathbb{R}\}$, that is, instead of a probability measure, we use the *Lebesgue measure*. Then,

$$D(f \| g) = \int_{-\infty}^{\infty} f(x) \log f(x) dx = -h(f).$$

Hence, for a probability law that contains a discrete component, its differential entropy is $-\infty$.

Mutual information

The mutual information between two jointly distributed r.v.'s X and Y can be defined via information divergence when they are not discrete:

- If (X, Y) is jointly distributed with PDF $f_{X,Y}$, then

$$\begin{aligned} I(X; Y) &= D(f_{X,Y} \| f_X \times f_Y) \\ &= E \left[\log \frac{f_{X,Y}(X, Y)}{f_X(X) f_Y(Y)} \right] \\ &= h(X) + h(Y) - h(X, Y) = h(X) - h(X|Y) = h(Y) - h(Y|X). \end{aligned}$$

- More generally, if $Y|X = x$ has a density for all x , then

$$I(X; Y) = D(f_{Y|X} \| f_Y | P_X) = h(Y) - h(Y|X).$$

Properties that extend to continuous r.v.'s

Proposition 1 (Chain rule)

Suppose (X, Y) is jointly distributed with joint PDF $f_{X,Y}$ and $h(X, Y)$ exists, then $h(X, Y) = h(X) + h(Y|X)$. More generally,

$$h(X^n) = \sum_{i=1}^n h(X_i | X^{i-1}).$$

Proposition 2 (Conditioning reduces differential entropy)

$$h(X|Y) \leq h(X), \quad h(X|Y, Z) \leq h(X|Z).$$

New properties of differential entropy

Differential entropy can be negative.

Since $b - a$ can be made arbitrarily small, $h(X) = \log(b - a)$ can be negative. Hence, the non-negative property of entropy *cannot* be extended to differential entropy.

Scaling will change the differential entropy.

Consider $X \sim \text{Unif}[0, 1]$. Then, $2X \sim \text{Unif}[0, 2]$. Hence,

$$h(X) = \log 1 = 0, \quad h(2X) = \log 2 = 1 \implies h(X) \neq h(2X).$$

This is in sharp contrast to entropy: $H(X) = H(g(X))$ as long as $g(\cdot)$ is an invertible function.

Scaling and translation

Proposition 3 (Scaling and Translation in the Scaler Case)

Let X be a continuous random variable with differential entropy $h(X)$.

- *Translation does not change the differential entropy: For a constant c , $h(X + c) = h(X)$.*
- *Scaling shifts the differential entropy: For a constant $a \neq 0$, $h(aX) = h(X) + \log |a|$.*

Proposition 4 (Scaling and Translation in the Vector Case)

Let \mathbf{X} be a continuous random vector with differential entropy $h(\mathbf{X})$.

- *For a constant vector \mathbf{c} , $h(\mathbf{X} + \mathbf{c}) = h(\mathbf{X})$.*
- *For an invertible matrix $\mathbf{a} \in \mathbb{R}^{n \times n}$, $h(\mathbf{a}\mathbf{X}) = h(\mathbf{X}) + \log |\det(\mathbf{a})|$.*

The proof of these propositions are left as exercises (simple calculus).

Differential entropy of Gaussian random vectors

Example 3 (Differential entropy of Gaussian random vectors)

For a n -dim random vector $\mathbf{X} \sim \mathcal{N}(\mathbf{m}, \mathbf{k})$, its differential entropy $h(\mathbf{X}) = \frac{1}{2} \log((2\pi e)^n \det(\mathbf{k}))$.

sol: For an n -dim random vector $\mathbf{X} \sim \mathcal{N}(\mathbf{m}, \mathbf{k})$, we can rewrite \mathbf{X} as

$$\mathbf{X} = \mathbf{a}\mathbf{W} + \mathbf{m},$$

where $\mathbf{a}\mathbf{a}^\top = \mathbf{k}$ and \mathbf{W} consists of $W_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$, $i = 1, \dots, n$.

Hence, by the translation and scaling properties of differential entropy:

$$\begin{aligned} h(\mathbf{X}) &= h(\mathbf{W}) + \log|\det(\mathbf{a})| = \sum_{i=1}^n h(W_i) + \frac{1}{2} \log \det(\mathbf{k}) \\ &= \frac{n}{2} \log(2\pi e) + \frac{1}{2} \log \det(\mathbf{k}) = \frac{1}{2} \log((2\pi e)^n \det(\mathbf{k})). \end{aligned}$$

Maximum differential entropy

Uniform distribution maximizes entropy for r.v. with finite support.

For differential entropy, the maximization problem needs to be associated with constraints on the distribution. (otherwise, it is simple to make it infinite)

It turns out that, under a **second moment constraint**, zero-mean Gaussian maximizes the differential entropy.

Theorem 1 (Maximum Differential Entropy under Covariance Constraint)

Let \mathbf{X} be a random vector with mean \mathbf{m} and covariance matrix

$$\mathbb{E}[(\mathbf{X} - \mathbf{m})(\mathbf{X} - \mathbf{m})^T] = \mathbf{k},$$

and \mathbf{X}^G be Gaussian with the same covariance \mathbf{k} . Then,

$$h(\mathbf{X}) \leq h(\mathbf{X}^G) = \frac{1}{2} \log((2\pi e)^n \det(\mathbf{k})).$$

pf: First, we can assume WLOG that both \mathbf{X} and \mathbf{X}^G are zero-mean, since translation does not change the differential entropy.

Let the PDF of \mathbf{X} be $f(x)$ and the PDF of \mathbf{X}^G be $f^G(x)$. Hence,

$$0 \leq D(f \| f^G) = E[\log f(\mathbf{X})] - E[\log f^G(\mathbf{X})] = -h(\mathbf{X}) - E_{\mathbf{X} \sim f}[\log f^G(\mathbf{X})].$$

Note that $\log f^G(x)$ is a quadratic function of x . \mathbf{X} and \mathbf{X}^G have the same second moment. Hence,

$$\begin{aligned} E_{\mathbf{X} \sim f}[\log f^G(\mathbf{X})] &= E_{\mathbf{X} \sim f^G}[\log f^G(\mathbf{X})] = -h(\mathbf{X}^G) \\ \implies 0 \leq D(f \| f^G) &= -h(\mathbf{X}) + h(\mathbf{X}^G) \\ \implies h(\mathbf{X}) &\leq h(\mathbf{X}^G). \end{aligned}$$

□

In general, maximum (differential) entropy problems subject to moment constraints can be solved via the non-negativity of KL divergence.

Summary

- Mutual information between two continuous r.v.'s X and Y with joint density $f_{X,Y}$: $I(X;Y) = \mathbb{E} \left[\log \frac{f_{X,Y}(X,Y)}{f_X(X)f_Y(Y)} \right]$.
- Differential entropy and conditional differential entropy:
 $h(X) := \mathbb{E} \left[\log \frac{1}{f_X(X)} \right]$, $h(X|Y) := \mathbb{E} \left[\log \frac{1}{f_{X|Y}(X|Y)} \right]$.
- $I(X;Y) = h(X) - h(X|Y) = h(Y) - h(Y|X)$.
- Chain rule, conditioning reduces differential entropy, non-negativity of mutual information and KL divergence: remain to hold.
- Differential entropy can be negative; $h(X) \not\leq h(X,Y)$.
- Translating a r.v. does not change the differential entropy, while scaling it shifts its differential entropy.
- Under a second moment constraint, zero-mean Gaussian maximizes the differential entropy.