

INFORMATION THEORY & CODING

Differential Entropy 2

Dr. Rui Wang

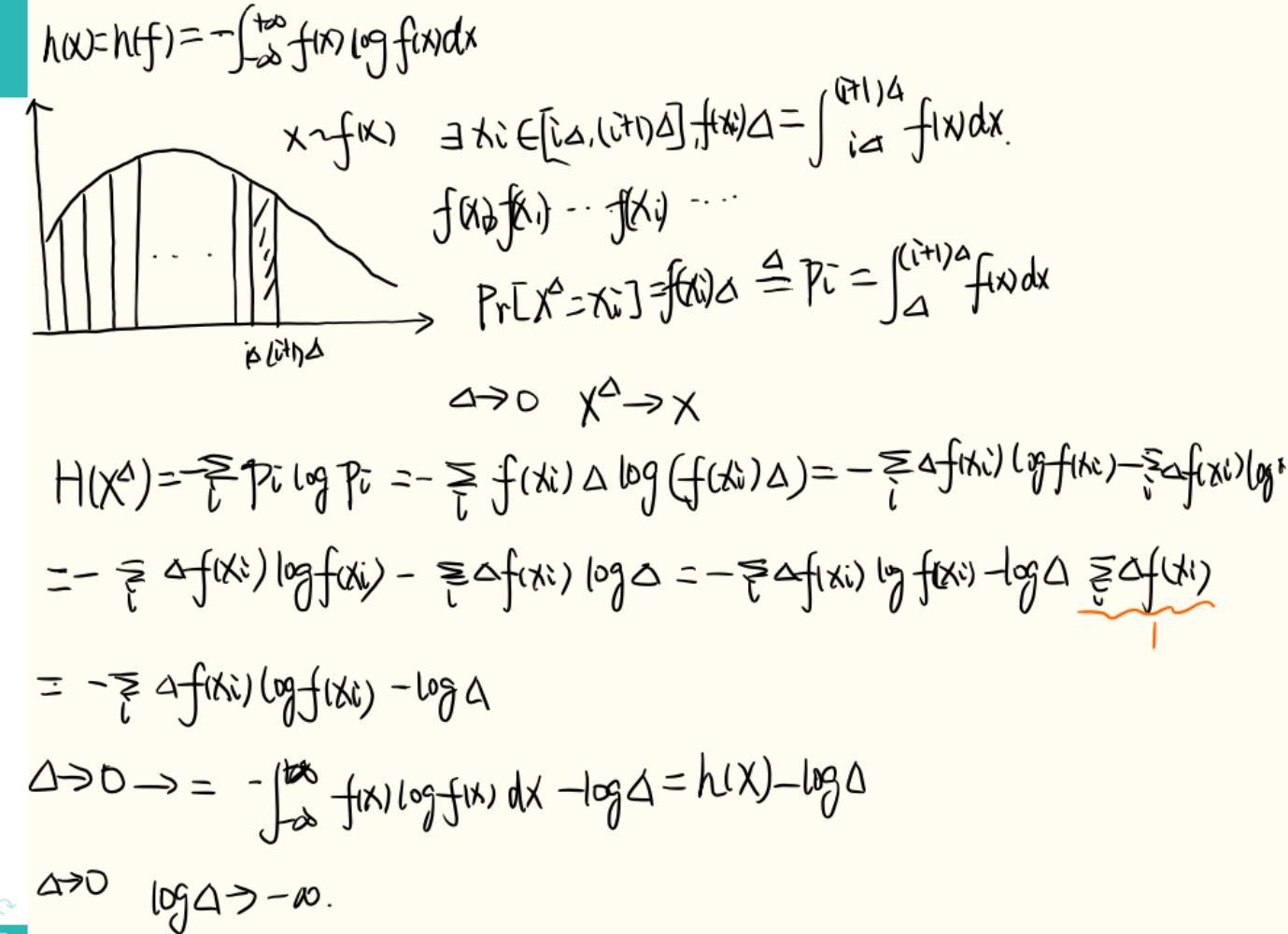
Department of Electrical and Electronic Engineering
Southern Univ. of Science and Technology (SUSTech)

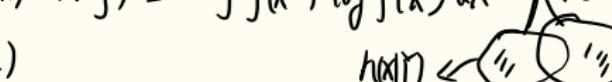
Email: wang.r@sustech.edu.cn

December 27, 2022



- Definitions
- AEP for Continuous Random Variables
- Relation of differential entropy to discrete entropy
- Joint and Conditional Differential Entropy
- Relative Entropy and Mutual Information
- Estimation Counterpart of Fano's Inequality



$$X_1, \dots, X_n \sim h(X_1, \dots, X_n) = h(f) = - \int f(x^n) \log f(x^n) dx^n$$


$$h(X|Y) = - \int f(x|y) \log f(x|y) dx dy = h(X, Y) - h(Y)$$

$$X_1, \dots, X_n \sim N(0, k) \quad k = E \left(\begin{matrix} X_1 \\ \vdots \\ X_n \end{matrix} \right) (X_1, \dots, X_n)$$

$$\phi_k(x_1, \dots, x_n) = \frac{1}{(\sqrt{2\pi})^n |K|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\vec{x} - \vec{u})^T K^{-1} (\vec{x} - \vec{u}) \right\} \quad \vec{x} = (x_1, \dots, x_n)$$

$$h(X_1, \dots, X_n) = - \int \phi(x^n) \log \phi(x^n) dx^n = - \int \phi(x^n) \left[\log(\sqrt{2\pi})^n |K|^{\frac{1}{2}} - \frac{1}{2} \log(\vec{x} - \vec{u})^T K^{-1} (\vec{x} - \vec{u}) \right] dx^n$$

$$= \frac{1}{2} \log[(2\pi)^n |K|] + \frac{\log e}{2} \int \phi(x^n) (\vec{x} - \vec{u})^T K^{-1} (\vec{x} - \vec{u}) dx^n$$

$$= \frac{1}{2} \log[(2\pi)^n |K|] + \frac{\log e}{2} E(\vec{x} - \vec{u})^T K^{-1} (\vec{x} - \vec{u})$$

$$= \frac{1}{2} \log[(2\pi)^n |K|] + \frac{\log e}{2} E \left\{ \text{tr} (\vec{x} - \vec{u})(\vec{x} - \vec{u})^T K^{-1} \right\}$$



Entropy of a multivariate Gaussian

Definition (Multivariate Gaussian Distribution)

If the joint pdf of X_1, X_2, \dots, X_n satisfies

$$f(\mathbf{x}) = f(x_1, \dots, x_n) = \frac{1}{(\sqrt{2\pi})^n |K|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T K^{-1}(\mathbf{x} - \mu)\right),$$

then X_1, X_2, \dots, X_n are multivariate/joint Gaussian/normal distributed with mean μ and covariance matrix K . Denote as $(X_1, X_2, \dots, X_n) \sim \mathcal{N}_n(\mu, K)$.

Theorem (Entropy of a multivariate normal distribution)

Let X_1, X_2, \dots, X_n have multivariate normal distribution with mean μ and covariance matrix K . Then

$$h(X_1, X_2, \dots, X_n) = h(\mathcal{N}_n(\mu, K)) = \frac{1}{2} \log(2\pi e)^n |K| \text{ bits},$$

where $|K|$ denotes the determinant of K .

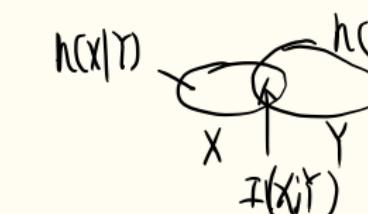
$$= \frac{1}{2} \log[(2\pi)^n |K|] + \log \text{tr} \tilde{E}(\tilde{x} - \mu)(\tilde{x} - \mu)^T K^{-1}$$

$$= \frac{1}{2} \log[(2\pi)^n |K|] + \frac{1}{2} \log e^n = \frac{1}{2} \log[(2\pi e)^n |K|]$$

$$f(x^n) g(x^n)$$

$$\text{Relative Entropy: } D(f||g) = \int f(x^n) \log \frac{f(x^n)}{g(x^n)} dx^n$$

$$\text{mutual Info. } I(X;Y) = \int f(x,y) \log \frac{f(x,y)}{f(x)f(y)} dx dy = D(f||g) + I(f(y)|f(x))$$



Relative entropy and mutual information

Definition

The relative entropy $D(f||g)$ between two pdfs f and g is

$$D(f||g) = \int f \log \frac{f}{g}.$$

Note: $D(f||g)$ is finite **only if** the support set of f is contained in the support set of g .

Definition

The mutual information $I(X; Y)$ between two random variables with joint pdf $f(x, y)$ is

$$I(X; Y) = \int f(x, y) \log \frac{f(x, y)}{f(x)f(y)} dx dy.$$

Relative entropy and mutual information

Definition

The **relative entropy** $D(f||g)$ between two pdfs f and g is

$$D(f||g) = \int f \log \frac{f}{g}.$$

Note: $D(f||g)$ is finite **only if** the support set of f is contained in the support set of g .

Definition

The **mutual information** $I(X; Y)$ between two random variables with joint pdf $f(x, y)$ is

$$I(X; Y) = \int f(x, y) \log \frac{f(x, y)}{f(x)f(y)} dx dy.$$

Relative entropy and mutual information

By definition, it is clear that

$$I(X;Y) = h(X) - h(X|Y) = h(Y) - h(Y|X) = h(X) + h(Y) - h(X,Y).$$

and

$$I(X;Y) = D\left(f(x,y) \middle\| f(x)f(y)\right).$$

Mutual information between correlated Gaussian r.v.s

- Let $(X, Y) \sim \mathcal{N}(0, K)$, where

$$K = \begin{bmatrix} \sigma^2 & \rho\sigma^2 \\ \rho\sigma^2 & \sigma^2 \end{bmatrix}.$$

- $h(X) = h(Y) = \frac{1}{2} \log(2\pi e) \sigma^2$
- $h(X, Y) = \frac{1}{2} \log(2\pi e)^2 |K| = \frac{1}{2} (\log 2\pi e)^2 \sigma^4 (1 - \rho^2)$
- $I(X; Y) = h(X) + h(Y) - h(X, Y) = -\frac{1}{2} \log(1 - \rho^2)$

if $\rho = 0$, X and Y are **independent**, the mutual information is 0.

if $\rho \pm 1$, X and Y are **perfectly correlated**, the mutual information is infinite.



- Let $(X, Y) \sim \mathcal{N}(0, K)$, where

$$K = \begin{bmatrix} \sigma^2 & \rho\sigma^2 \\ \rho\sigma^2 & \sigma^2 \end{bmatrix}.$$

- $h(X) = h(Y) = \frac{1}{2} \log(2\pi e) \sigma^2$
- $h(X, Y) = \frac{1}{2} \log(2\pi e)^2 |K| = \frac{1}{2} (\log 2\pi e)^2 \sigma^4 (1 - \rho^2)$
- $I(X; Y) = h(X) + h(Y) - h(X, Y) = -\frac{1}{2} \log(1 - \rho^2)$

if $\rho = 0$, X and Y are **independent**, the mutual information is 0.

if $\rho \pm 1$, X and Y are **perfectly correlated**, the mutual information is infinite.

Theorem

$D(f||g) \geq 0$ with *equality* iff $f = g$ almost everywhere.

Proof.

Let \mathcal{S} be the support set of f . Then

$$\begin{aligned} -D(f||g) &= \int_{\mathcal{S}} f \log \frac{g}{f} \\ &\leq \log \int_{\mathcal{S}} f \frac{g}{f} \quad (\text{by Jensen's inequality}) \\ &= \log \int_{\mathcal{S}} g \\ &\leq \log 1 = 0 \end{aligned}$$

Theorem

$D(f||g) \geq 0$ with *equality* iff $f = g$ almost everywhere.

Proof.

Let \mathcal{S} be the support set of f . Then

$$\begin{aligned} -D(f||g) &= \int_{\mathcal{S}} f \log \frac{g}{f} \\ &\leq \log \int_{\mathcal{S}} f \frac{g}{f} \quad (\text{by Jensen's inequality}) \\ &= \log \int_{\mathcal{S}} g \\ &\leq \log 1 = 0 \end{aligned}$$



Properties of differential entropy

- $I(X;Y) \geq 0$ with **equality** iff X and Y are independent.
- $h(X|Y) \leq h(X)$ with **equality** iff X and Y are independent.

Theorem (Chain rule for differential entropy)

$$h(X_1, X_2, \dots, X_n) = \sum_{i=1}^n h(X_i | X_1, X_2, \dots, X_{i-1}).$$

- $h(X_1, X_2, \dots, X_n) \leq \sum h(X_i)$, with **equality** iff X_1, X_2, \dots, X_n are independent.

Properties of differential entropy

Theorem (Translation does not change the differential entropy)

$$h(X + c) = h(X).$$

Theorem

$$h(aX) = h(X) + \log |a|.$$

Proof.

Let $Y = aX$, Then $f_Y(y) = \frac{1}{|a|} f_X(\frac{y}{a})$, and we have

$$\begin{aligned} h(aX) &= - \int f_Y(y) \log f_Y(y) dy = - \int \frac{1}{|a|} f_X\left(\frac{y}{a}\right) \log \left(\frac{1}{|a|} f_X\left(\frac{y}{a}\right)\right) dy \\ &= - \int f_X(x) \log f_X(x) dx + \log |a| = h(X) + \log |a| \end{aligned}$$

Properties of differential entropy

Theorem (Translation does not change the differential entropy)

$$h(X + c) = h(X).$$

Theorem

$$h(aX) = h(X) + \log |a|.$$

Proof.

Let $Y = aX$, Then $f_Y(y) = \frac{1}{|a|} f_X(\frac{y}{a})$, and we have

$$\begin{aligned} h(aX) &= - \int f_Y(y) \log f_Y(y) dy = - \int \frac{1}{|a|} f_X(\frac{y}{a}) \log \left(\frac{1}{|a|} f_X \left(\frac{y}{a} \right) \right) dy \\ &= - \int f_X(x) \log f_X(x) dx + \log |a| = h(X) + \log |a| \end{aligned}$$

Theorem (Translation does not change the differential entropy)

$$h(X + c) = h(X).$$

Theorem

$$h(aX) = h(X) + \log |a|.$$

Corollary.

$$h(A\mathbf{X}) = h(\mathbf{X}) + \log |\det(A)|.$$

□



期望
协方差矩阵

$$0 \leq D(g || \phi_k) = \int g(x^n) \log \frac{g(x^n)}{\phi(x^n)} dx^n = -h(g) - \int g(x^n) \log \phi(x^n) dx^n$$

$$\int g(x^n) \log \phi(x^n) dx^n = \int g(x^n) \left[\log \frac{1}{(\sqrt{2\pi})^n |K|^{\frac{1}{2}}} - \frac{1}{2} \log e^{\frac{1}{2}(\vec{x}-\vec{\mu})^T K^{-1}(\vec{x}-\vec{\mu})} \right] dx^n$$

$$= \int g(x^n) \left[\log \frac{1}{(\sqrt{2\pi})^n |K|^{\frac{1}{2}}} \right] dx^n - \frac{1}{2} \log \underbrace{\int g(x^n) (\vec{x}-\vec{\mu})^T K^{-1} (\vec{x}-\vec{\mu}) dx^n}_{E_g[(\vec{x}-\vec{\mu})^T K^{-1} (\vec{x}-\vec{\mu})]}$$

$$= E_g[\text{tr}(\vec{x}-\vec{\mu})(\vec{x}-\vec{\mu})^T K^{-1}]$$

$$= \text{tr} K^{-1} = \text{tr} I = E[(\vec{x}-\vec{\mu})^T K^{-1} (\vec{x}-\vec{\mu})]$$

Theorem

Let the random vector $\mathbf{X} \in \mathbb{R}^n$ have zero mean and covariance $K = \mathbb{E} \mathbf{X} \mathbf{X}^T$ (i.e., $K_{ij} = \mathbb{E} X_i X_j$, $1 \leq i, j \leq n$). Then

$$h(\mathbf{X}) \leq \frac{1}{2} \log(2\pi e)^n |K|$$

with **equality** iff $\mathbf{X} \sim \mathcal{N}(0, K)$.

AWGN.

$$X \xrightarrow{\text{AWGN}} Y = X + Z \sim N(0, \sigma^2)$$

$$\max_{f(x)} I(X; Y) = C \quad I(X; Y) = h(Y) - h(Y|X) = h(Y) - h(Z|X) = h(Y) - h(Z)$$

$$h(Z) = \frac{1}{2} \log 2\pi e \sigma^2 \quad \text{因为 } E[Z] = 0$$

$$E[Y^2] = E[(X+Z)^2] = E[X^2 + 2XZ + Z^2] = E[X^2] + E[Z^2] = p + \sigma^2$$

$$\Rightarrow h(Y) = \frac{1}{2} \log 2\pi e (p + \sigma^2)$$

$$\Rightarrow I(X; Y) = h(Y) - h(Z) = \frac{1}{2} \log (p + \frac{p}{\sigma^2})$$

Random variable X , estimator \hat{X} . The expected prediction error $E(X - \hat{X})^2$.

Theorem (Estimation error and differential entropy)

For any random variable X and estimator \hat{X} ,

$$\mathbb{E}(X - \hat{X})^2 \geq \frac{1}{2\pi e} \exp(2h(X)),$$

with *equality* iff X is Gaussian and \hat{X} is the *mean* of X .

Theorem (Estimation error and differential entropy)

For any random variable X and estimator \hat{X} ,

$$\mathbb{E}(X - \hat{X})^2 \geq \frac{1}{2\pi e} \exp\left(2h(X)\right),$$

with *equality* iff X is Gaussian and \hat{X} is the *mean* of X .

Proof.

We have

$$\begin{aligned}\mathbb{E}(X - \hat{X})^2 &\geq \min_{\hat{X}} \mathbb{E}(X - \hat{X})^2 \\ &= \mathbb{E}(X - \mathbb{E}(X))^2 \quad \text{mean is the best estimator} \\ &= \text{Var}(X) \\ &\geq \frac{1}{2\pi e} \exp\left(2h(X)\right). \quad \text{The Gaussian has maximum entropy}\end{aligned}$$

Summary

- Discrete r.v. \Rightarrow continuous r.v.
- entropy \Rightarrow differential entropy.
- Many things similar: mutual information, relative entropy, AEP, chain rule, ...

Some things different: $h(X)$ can be negative, maximum entropy distribution is Gaussian

