So, if $R > C$ then $P_e^{(n)}$ is bounded away from zero for sufficiently large $n$, and, hence for all $n$. The latter comes from the fact that if $P_e^{(n)} \to 0$ for small $n$, we could construct long codes by concatenating short codes.

X Lecture 7, Oct. 14, 2003

## Continuous Sources

A continuous random variable $X$ is characterised by a Cumulative Distribution Function $F_X(x)$,

$$F_X(x) = P_r(X \le x)$$

the derivative of $F_X(x)$ w.r.t. $x$ is called the <u>probability density function</u> (pdf),

$$f_X(x) = \frac{\partial F_X(x)}{\partial x}$$

when there is no ambiguity, we drop the subscript and denote the pdf by $f(x)$.

It is clear that,

$$F(x) = \int_{-\infty}^{x} f(x) \, dx$$

therefore

$$\int_S f(x) \, dx = 1$$

<u>83</u>

where $S$ is the support of the random variable $X$, i.e., the set of $x$ such that $f(x) > 0$.

Assume that a source generates $X_1, X_2, \ldots$ i.i.d. and $\sim f(x)$. Then, for any $n$-tuple $(x_1, \ldots x_n) = x^n$

$$f(x_1, x_2, \ldots, x_n) = \prod_{i=1}^{n} f(x_i)$$

Taking logarithm of $f(x_1, \ldots x_n)$ we have,

$$\log f(x_1, \ldots, x_n) = \sum_{i=1}^{n} \log f(x_i)$$

normalizing by $n$, we have

$$-\frac{1}{n} \log f(x_1, \ldots x_n) = \frac{-1}{n} \sum_{i=1}^{n} \log f(x_i)$$

The right-hand side of the above equation is the sample mean of $\log f(x)$ and according to law of large numbers, it tends to the expectation of $\log f(x)$ as $n \to \infty$. So,

$$-\frac{1}{n} \log f(x_1, \ldots x_n) \xrightarrow[n \to \infty]{} -E[f(x)] = -\int_S f(x) \log f(x) \, dx$$
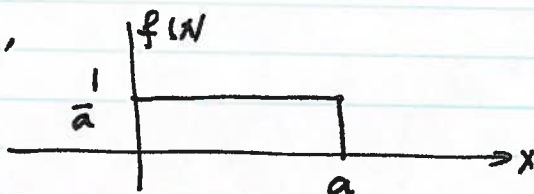
We denote this expectation by $h(X)$ or $h(f)$ and call $h(X)$ differential entropy:

$$h(X) = -E[\log f(x)] = -\int_S f(x) \log f(x) \, dx$$

$h(X)$ is similar to $H(X)$ for discrete alphabet sources.

however, fails to have some of the properties of $H(X)$ such as positivity which make $H(X)$ a measure of information.

Example : differential entropy for a uniformly distributed source.



$$h(X) = -\int_0^a \frac{1}{a} \log \frac{1}{a} \, dx = \log a$$

Note that for $a > 1$ as $a$ increases so does $h(X)$ and this is in agreement with $h(X)$ being a measure of uncertainty as increasing $a$ makes the particular value of $X$ more uncertain and, therefore, more information-beari

However, we also note that for $a < 1$, $h(X)$ is negative On the other hand, $2^{h(X)} = 2^{\log a} = a$ is the volume of the support set of $X$ which is always non-negative.

Example : $h(X)$ for a Gaussian source.

Here $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/2\sigma^2}$

Then

$$h(X) = -\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/2\sigma^2} \log \left[ \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/2\sigma^2} \right] dx$$

85

or,

$$h(X) = -\int_{-\infty}^{\infty} f(x) \log\left[\frac{1}{\sqrt{2\pi}\,\sigma}\right] dx$$

$$+ \int_{-\infty}^{\infty} f(x) \left[\frac{x^2}{2\sigma^2}\right] \log_2 e \, dx$$

$$= \frac{1}{2}\log(2\pi\sigma^2) + \frac{\log_2 e \int_{-\infty}^{\infty} f(x) x^2 dx}{2\sigma^2}$$

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

$$= \frac{1}{2}\log(2\pi e \sigma^2) \qquad \text{bits}$$

Definition : For any $\varepsilon > 0$ and any $n$, define the typical set $A_\varepsilon^{(n)}$ w.r.t. $f(x)$ as:

$$A_\varepsilon^{(n)} = \left\{ (x_1, x_2, \ldots, x_n) \in S^n : \left| -\frac{1}{n}\log f(x_1, \ldots x_n) - h(X)\right| \leq \varepsilon \right\}$$

Theorem (AEP): The typical set $A_\varepsilon^{(n)}$ has the following properties:

1) $P_r(A_\varepsilon^{(n)}) \geq 1 - \varepsilon$ for $n$ sufficiently large

2) $Vol(A_\varepsilon^{(n)}) \leq 2^{n(h(X)+\varepsilon)}$ for all $n$

3) $Vol(A_\varepsilon^{(n)}) \geq (1-\varepsilon) 2^{n(h(X)-\varepsilon)}$ for $n$ sufficiently large.

where $Vol(A)$ for a set $A \in R^n$ is defined as

$$Vol(A) = \int_A dx_1 \, dx_2 \cdots dx_n.$$

Proof:

$$-\frac{1}{n} \log f(x_1, x_2, \ldots, x_n) = -\frac{1}{n} \sum_i f(x_i) \to h(x)$$

So, for any $\varepsilon > 0$ there is some $n_o$ such that for any $n > n_o$, we have

$$Pr \left\{ | -\frac{1}{n} \log f(x_1, \ldots, x_n) - h(x) | \not> \varepsilon \right\} \leq \varepsilon$$

or

$$Pr \left\{ (x_1, \ldots, x_n) \notin A_\varepsilon^{(n)} \right\} \leq \varepsilon$$

and, hence,

$$Pr \left\{ (x_1, \ldots x_n) \in A_\varepsilon^{(n)} \right\} > 1 - \varepsilon$$

This proves part 1.

For part 2:

$$1 = \int_S f(x_1, \ldots, x_n) \, dx_1 \, dx_2 \cdots dx_n$$

$$\geq \int_{A_\varepsilon^{(n)}} f(x_1, \ldots x_n) \, dx_1 \ldots dx_n$$

$$\geq \int_{A_\varepsilon^{(n)}} 2^{-n(h(X) + \varepsilon)} \, dx_1 \ldots dx_n$$

$$= 2^{-n(h(X) + \varepsilon)} Vol(A_\varepsilon^{(n)}) \qquad QED \text{ [Property 2]}$$

87

## Property 3:

if $n$ is large enough so that property 1 holds, then:

$$1 - \epsilon \leq \int_{A_\epsilon^{(n)}} f(x_1, \ldots x_n) \, dx_1 \cdots dx_n$$

$$\leq \int_{A_\epsilon^{(n)}} 2^{-n(h(X) - \epsilon)} \, dx_1 \cdots dx_n$$

$$= 2^{-n(h(X) - \epsilon)} \, \text{Vol}(A_\epsilon^{(n)}) \qquad \text{QED [Property 3]}$$

~~~~~~~~~

## Joint and Conditional differential entropy:

**Definition:** The differential entropy of a set $X_1, X_2, \ldots X_n$ of random variables is

$$h(X_1, \ldots X_n) = -\int_{S^n} f(\underline{x}^n) \log f(\underline{x}^n) \, d\underline{x}^n$$

**Definition:** If $X$ and $Y$ have joint density function $f(x, y)$, we can define the conditional differential entropy $h(X|Y)$ as,

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

substituting $f(x|y) = \dfrac{f(x, y)}{f(y)}$, we get

$$h(X|Y) = h(X, Y) - h(Y)$$

Example: The ~~entire~~ differential entropy of a multivariate normal sequence is

$$h(X_1, X_2, \ldots, X_n) = \frac{1}{2} \log (2\pi e)^n |K| \quad \text{bits}$$

where $|K|$ is the determinant of ~~a~~ the covariance matrix of $X_1, X_2, \ldots X_n$.

Proof:

$$f(\underline{x}^n) = \frac{1}{(\sqrt{2\pi})^n |K|^{1/2}} \exp\left[-\frac{1}{2}(\underline{x}^n - \underline{\mu}^n)^T K^{-1}(\underline{x}^n - \underline{\mu}^n)\right]$$

Then,

$$h(\underline{x}^n) = - \int f(\underline{x}^n) \log f(\underline{x}^n) \, d\underline{x}^n$$

$$= \frac{\log e}{2} \int f(\underline{x}^n) \left[(\underline{x}^n - \underline{\mu}^n) K^{-1}(\underline{x}^n - \underline{\mu}^n)\right] d\underline{x}^n$$

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

$$= \frac{\log e}{2} E\left[(\underline{x}^n - \underline{\mu}^n) K^{-1}(\underline{x}^n - \underline{\mu}^n)\right] + \frac{1}{2} \log(2\pi)^n |K|$$

$$= \frac{\log e}{2} E\left[\sum_i \sum_j (x_i - \mu_i) K^{-1}_{ij} (x_j - \mu_j)\right] + \frac{1}{2} \log (2\pi)^n |K|$$

$$= \frac{\log e}{2} \sum_i \sum_j E\left[(x_i - \mu_i)(x_j - \mu_j)\right] K^{-1}_{ij} + \frac{1}{2} \log (2\pi)^n |K|$$

$$= \frac{\log e}{2} \sum_i \sum_j K_{ij} K^{-1}_{ji} + \frac{1}{2} \log(2\pi)^n |K|$$

$$= \frac{\log e}{2} \sum_i (K^{-1}K)_{ii} + \frac{1}{2} \log (2\pi)^n |K|$$

$$= \frac{\log e}{2} \sum_i I_{ii} + \frac{1}{2} \log (2\pi)^n |K|$$

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

$$= \frac{1}{2} \log (2\pi e)^n |K| \qquad\qquad QED$$

## Mutual information:

<u>Definition</u> : The mutual information $I(X;Y)$ between $X$ and $Y \sim f(x,y)$ is defined as:

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

$$= \int_S f(x,y) \log \frac{f(y|x)}{f(y)} \, dx\, dy$$

It is easy to show that

$$I(X;Y) = h(X) - h(X|Y) = h(Y) - h(Y|X)$$

Note that

$$I(X;Y) = D(f(x,y) \| f(x)f(y))$$

The properties of mutual information and diff. entr

1) $I(X;Y) \geq 0$ with equality iff $X$ and $Y$ are independent.

2) $h(X|Y) \leq h(X)$ with equality if $X$ and $Y$ are independent.

3) The chain rule for differential entropy
$$h(X_1, X_2, \ldots, X_n) = \sum_{i=1}^{n} h(X_i | X_1, X_2, \ldots, X_{i-1})$$

4) $h(X_1, X_2, \ldots, X_n) \leq \sum_{i=1}^{n} h(X_i)$
with equality if $X_1, \ldots X_n$ are independent.

<u>Hadamard inequality</u>: Let $X^n \sim N(0,K)$ be a multivariate normal random variable, then
$$|K| \leq \prod_{i=1}^{n} K_{ii}$$

Proof:
$$h(X_1, \ldots, X_n) = \frac{1}{2} \log(2\pi e)^n |K|$$
$$h(X_i) = \frac{1}{2} \log(2\pi e K_{ii})$$
$$\frac{1}{2} \log(2\pi e)^n |K| \leq \frac{1}{2} \sum_{i=1}^{n} \log(2\pi e K_{ii})$$
$$\Rightarrow \log(2\pi e)^n |K| \leq \log \prod_{i=1}^{n} (2\pi e K_{ii})$$

$$\log (2\pi e)^n |K| \leq \log (2\pi e)^n \prod_{i=1}^{n} K_{ii}$$

$$\Rightarrow \quad |K| \leq \prod_{i=1}^{n} K_{ii}$$

an example situation where this inequality is useful is the MIMO channels where correlation between different paths results in reduction in capacity.

**Theorem** $\quad h(X+c) = h(X)$

with suppor

Proof: if $X \sim f(x)$ then $X+c \sim f(x+c)$ ~~over~~ $S+c$

$$h(X+c) = - \int_{S+c} f(x+c) \log f(x+c) \, dx = -\int_{S} f(x) \log f(x+c)$$

$$= h(X).$$

**Theorem:** $\quad h(aX) = h(X) + \log |a|$

Let $Y = aX$. Then $f_Y(y) = \frac{1}{|a|} f_X\left(\frac{y}{a}\right)$

$$h(aX) = - \int f_Y(y) \log f_Y(y) \, dy$$

$$= - \int \frac{1}{|a|} f_X\left(\frac{y}{a}\right) \log\left(\frac{1}{|a|} f_X\left(\frac{y}{a}\right)\right) dy$$

$$= - \int f_X(x) \log f_X(x) \, dx + \log |a|$$

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

<u>Corollary</u> : for a random vector $\underline{X}$, we have

$$h(\underline{A}\,\underline{X}) = h(\underline{X}) + \log|A|$$

where $|A|$ is the determinant of the matrix $A$.

The most difficult source to ~~include~~ Compress:
The following theorem indicates that among all ~~the~~ random vectors with zero mean and Common covariance matrix $K = E[\underline{X}\,\underline{X}^T]$, the Gaussian vector has the largest differential entropy. This implies that the most difficult source to Compress is a Gaussian Source.

<u>Theorem</u> : Let the random vector $X \in \mathbb{R}^n$ have zero mean and Covariance $K = E[XX^T]$, i.e.,

$K_{ij} = E[X_i X_j]$, $i,j = 1, 2, \ldots, n$. Then

$$h(\underline{X}) \leq \tfrac{1}{2} \log(2\pi e)^n |K|.$$

<u>Proof</u>:

Let $p(\underline{x})$ be the pdf of $\underline{X}$
and $q(\underline{x}) = \dfrac{1}{(\sqrt{2\pi}\,)^n |K|^{1/2}} e^{-\frac{1}{2}\underline{x}^T K^{-1}\underline{x}}$

93

Then,

$$\int_{\mathbb{R}^n} p(\underline{x}) \log \frac{p(\underline{x})}{q(\underline{x})} \, d\underline{x} = \int_{\mathbb{R}^n} p(\underline{x}) \ln\left[\frac{p(\underline{x})}{q(\underline{x})}\right] \log e \, d\underline{x}$$

$$\geq \int_{\mathbb{R}^n} p(\underline{x})\left[1 - \frac{q(\underline{x})}{p(\underline{x})}\right] \log e \, d\underline{x} = 0$$

But,

$$\int_{\mathbb{R}^n} p(\underline{x}) \log \frac{p(\underline{x})}{q(\underline{x})} \, dx = \int_{\mathbb{R}^n} p(\underline{x}) \log p(\underline{x}) \, dx$$

$$- \int_{\mathbb{R}^n} p(\underline{x}) \log q(\underline{x}) \, d\underline{x} = -h(x) - \int_{\mathbb{R}^n} p(x) \log \frac{1}{\sqrt{2\pi}^n |K|^{1/2}} \, dx$$

$$+ \frac{\log e}{2} \int_{\mathbb{R}^n} (x^T K^{-1} \underline{x}) \, p(\underline{x}) \, d\underline{x}$$

$$= -h(\underline{x}) + \frac{1}{2} \log (2\pi)^n |K| + \frac{\log e}{2} \int_{\mathbb{R}^n} \left(\sum_i \sum_j x_i x_j (K^{-1})_{ij}\right) p(\underline{x})$$

$$= -h(\underline{x}) + \frac{1}{2} \log (2\pi)^n |K| + \frac{\log e}{2} \sum_i \sum_j \left[\int_{\mathbb{R}^n} x_i x_j \, p(\underline{x}) \, d\underline{x}\right] (K)_{ij}^{-1}$$

$$= -h(\underline{x}) + \frac{1}{2} \log (2\pi)^n |K| + \frac{\log e}{2} \sum_i \sum_j K_{ji} (K^{-1})_{ij}$$

$$= -h(\underline{x}) + \frac{1}{2} \log (2\pi)^n |K| + \frac{\log e}{2} \underbrace{\sum_i (K^{-1}K)_{ii}}_{n}$$

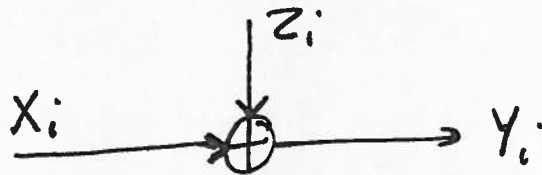$$= -h(\underline{x}) + \frac{1}{2} \log (2\pi e)^n |K| \geq 0$$

94

so,

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

X Lecture 8, Oct. 21, 2003

## Gaussian Channel

The channel we consider here is a discrete time channel with inputs $X_1, X_2, \ldots$ and output $Y_1, Y_2, \ldots$ where

$$Y_i = X_i + Z_i$$

where $Z_i \sim N(0, N)$



if the noise power (variance) $N$ is zero or the transmission power is limitless, Then it is possible to transmit an infinite number of bits per use. In such a case (unconstrained) capacity of the channel is infinite. In practice, however, there is a limit on the transmi