

- Lecture 2:

Markov Chain

Random variables X, Y, Z form a ~~an~~ Markov Chain if the conditional distribution of Z depends only on Y and is conditionally independent of X . That is if one knows the value of Y then X cannot say anything more about Z .

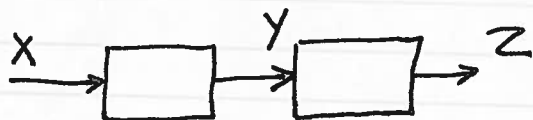
We show this as $X \rightarrow Y \rightarrow Z$ and formally, we have $X \rightarrow Y \rightarrow Z$ if

$$p(x, y, z) = p(x) p(y|x) p(z|y)$$

if $X \rightarrow Y \rightarrow Z$ then X and Z are conditionally independent, i.e.,

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

Data Processing Inequality



If $X \rightarrow Y \rightarrow Z$ then $I(X; Y) \geq I(X; Z)$

Proof:

$$I(X; Y, Z) = I(X; Z) + I(X; Y|Z) = I(X; Y) + I(X; Z|Y)$$

$I(X; Z|Y) = 0$ since X and Z are conditionally independent.

So,

$$I(X; Y) = I(X; Z) + I(X; Y|Z) \geq I(X; Z)$$

In particular,
if $Z = g(Y)$, then:

$$I(X; Y) \geq I(X; g(Y))$$

Since $X \rightarrow Y \rightarrow g(Y)$

Corollary: If $X \rightarrow Y \rightarrow Z$ then:

$$I(X; Y|Z) \leq I(X; Y)$$

Proof:

$$\begin{aligned} I(X; Y, Z) &= I(X; Y) + I(X; Z|Y) \\ &= I(X; Z) + I(X; Y|Z) \end{aligned}$$

$$I(X; Y|Z) = I(X; Y) - I(X; Z) \leq I(X; Y)$$

Sufficient Statistics

Assume that we have a ^{family} probability mass function $\{f_{\theta}(x)\}$ indexed by θ . Let X be a sample from a distribution in this family. Take $T(X)$ to be any statistic (function of the sample)

like sample mean, sample variance, etc.

Then $\theta \rightarrow X \rightarrow T(X)$ and by data processing theorem:

$$I(\theta; T(X)) \leq I(\theta; X)$$

A statistic $T(X)$ is called sufficient if the above holds with equality.

Definition: A function $T(X)$ is a sufficient statistic if X is independent of θ given $T(X)$, i.e.,

$\theta \rightarrow T(X) \rightarrow X$. This is equivalent to:

$$I(\theta; X) = I(\theta; T(X))$$

for all distributions on θ .

Example:

$$f_{\theta}(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\theta)^2}{2}} = N(\theta, \sigma)$$

Let X_1, \dots, X_n be drawn independently according to this distribution. Then the sufficient statistics

for θ is $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. It can be shown that the conditional distribution of X_1, \dots, X_n given \bar{X}_n and n does not depend on θ .

Fano's inequality

Fano's inequality relates the probability of error and Channel equivocation.

That is $P_e = P(\hat{X} \neq X)$ and $H(X|Y)$

Suppose we wish to estimate X based on the observation of Y . X is distributed as $p(x)$ and Y is related to X through $p(y|x)$, then:

Theorem:

$$H(P_e) + P_e \log(|X| - 1) \geq H(X|Y)$$

a little bit weaker version is

$$|H(P_e) + P_e \log |X| \geq H(X|Y)$$

or

$$P_e \geq \frac{H(X|Y) - 1}{\log |X|}$$

Proof:

Define r.v. $E = \begin{cases} 1 & \text{if } \hat{X} \neq X \\ 0 & \text{if } \hat{X} = X \end{cases}$

Using the chain rule, we have:

$$\begin{aligned} H(E, X|Y) &= H(X|Y) + H(E|X, Y) \\ &= H(E|Y) + H(X|E, Y) \end{aligned}$$

$$H(E|Y) \leq H(E) = H(P_e)$$

$$\begin{aligned} H(X|E, Y) &= P_r(E=0) H(X|Y, E=0) \\ &\quad + P_r(E=1) H(X|Y, E=1) \end{aligned}$$

$$\leq (1-P_e) 0 + P_e \log(|X|-1)$$

So,

$$H(X|Y) \leq H(P_e) + P_e \log(|X|-1)$$

The Asymptotic Equipartition Property (AEP)

AEP is the ^{application} ~~statement~~ of the law of large numbers in information theory. The same way that $\frac{1}{n} \sum_{i=1}^n X_i \rightarrow E[X]$ we show that

$$-\frac{1}{n} \log p(X_1, \dots, X_n) \rightarrow H(X)$$

This means that $p(X_1, \dots, X_n)$ will be close to $2^{-nH(X)}$.

Theorem (AEP): If X_1, \dots, X_n are i.i.d. $\sim p(x)$ then

$$-\frac{1}{n} \log p(X_1, \dots, X_n) \rightarrow H(X) \text{ in probability}$$

that is for any $\epsilon > 0$ there is some n_0 such ^{$\delta > 0$} that for any $n > n_0$

$$Pr \left[\left| -\frac{1}{n} \log p(X_1, \dots, X_n) - H(X) \right| > \epsilon \right] < \delta$$

Proof:

$$-\frac{1}{n} \log p(X_1, \dots, X_n) = -\frac{1}{n} \sum_{i=1}^n \log p(X_i)$$

$$\rightarrow -E[\log p(X_i)] \quad \leftarrow \text{(using weak law of large numbers)}$$

$$= H(X)$$

Definition The typical set $A_\epsilon^{(n)}$ (ϵ -typical set) with respect to $p(x)$ is defined as the set of all sequences $(x_1, x_2, \dots, x_n) \in X^n$ such that

$$2^{-n(H(X)+\epsilon)} \leq p(x_1, \dots, x_n) \leq 2^{-n(H(X)-\epsilon)}$$

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

1) for each $(x_1, \dots, x_n) \in A_\epsilon^{(n)}$

$$H(X) - \epsilon \leq -\frac{1}{n} \log p(x_1, \dots, x_n) \leq H(X) + \epsilon$$

2) $P\{A_\epsilon^{(n)}\} > 1 - \epsilon$ for sufficiently large n .

3) $|A_\epsilon^{(n)}| \leq 2^{n(H(X)+\epsilon)}$ where $|A|$ denotes the cardinality of A .

4) $|A_\epsilon^{(n)}| \geq (1 - \epsilon) 2^{n(H(X)-\epsilon)}$

Proof:

1) from the definition

$$-n(H(X)+\epsilon) \leq \log p(x_1, \dots, x_n) \leq -n(H(X)-\epsilon)$$

$$H(X) - \epsilon \leq -\frac{1}{n} \log p(x_1, \dots, x_n) \leq H(X) + \epsilon$$

2) For sufficiently large n

$$Pr\{A_\epsilon^{(n)}\} = Pr\{(x_1, \dots, x_n) \in A_\epsilon^{(n)}\}$$

$$= P\left\{ \left| -\frac{1}{n} \log p(x_1, \dots, x_n) - H(X) \right| \leq \epsilon \right\}$$

$$= 1 - P\left\{ \left| -\frac{1}{n} \log p(x_1, \dots, x_n) - H(X) \right| > \epsilon \right\} \geq 1 - \delta$$

by choosing $\delta = \epsilon$ ~~the~~ part 2 is proven.

3)

$$Pr\{A_\epsilon^{(n)}\} = \sum_{A_\epsilon^{(n)}} p(x) \leq \sum_{X^n} p(x) = 1$$

so

$$1 \geq \sum_{A_\epsilon^{(n)}} p(x) \geq \sum_{A_\epsilon^{(n)}} 2^{-n(H(X)+\epsilon)} = 2^{-n(H(X)+\epsilon)} |A_\epsilon^{(n)}|$$

and, therefore,

$$|A_\epsilon^{(n)}| \leq 2^{n(H(X)+\epsilon)}$$

4) $Pr(A_\epsilon^{(n)}) \geq 1 - \epsilon$ from property 2

so

$$1 - \epsilon \leq Pr(A_\epsilon^{(n)}) = \sum_{A_\epsilon^{(n)}} p(x) \leq \sum_{A_\epsilon^{(n)}} 2^{-n(H(X)-\epsilon)} = 2^{-n(H(X)-\epsilon)} |A_\epsilon^{(n)}|$$

Therefore:

$$|A_\epsilon^{(n)}| \geq (1 - \epsilon) 2^{-n(H(X)-\epsilon)}$$

in summary:

$$(1 - \epsilon) 2^{-n(H(X)-\epsilon)} \leq |A_\epsilon^{(n)}| \leq 2^{-n(H(X)+\epsilon)}$$

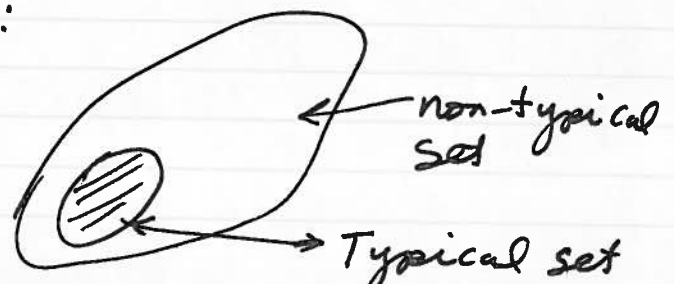
Interpretation:

realistic

The AEP indicates that the number of possibilities of a sequence of i.i.d. events is, usually, much less than the $|X|^n$ possibilities. It also indicates that all those possibilities that usually occur have roughly the same probability. In essence, it reveals the "dullness" of the world in the sense that ^{first of all,} out of a huge number of possibilities only ~~quite~~ just a few happen in practice and secondly, those occurring ~~at the~~ ^{roughly} cause the same level of "surprise".

Application of AEP in data Compression:

Assume that X_1, \dots, X_n are $n p(X)$, i.i.d. and we wish to ~~compress~~ ^{we divide} represent them. The total $|X|^n$ sequences ~~is~~ ^{is} divided into the set of typical sequences $A_\epsilon^{(n)}$ having less than $2^{n(H(X) + \epsilon)}$ members and the set of non-typical sequences.



To represent each typical sequence, we need at most $\log_2 2^{n(H(X)+\epsilon)} + 1 = n(H(X)+\epsilon) + 1$ bits. We devise a coding scheme as follows:

for each sequence, we specify the sequence to be a typical or non-typical by adding preamble bit, say: 0 for typical and 1 for non-typical

then, we use:

1) $n(H(X)+\epsilon) + 1$ bits for each typical sequence. We can use any arbitrary indexing for counting the $2^{n(H(X)+\epsilon)}$ typical sequences.

2) $n \log |X| + 1$ per non-typical sequence.

So:

$$\begin{aligned}
 E[l(X^n)] &= \sum_{x^n} p(x^n) l(x^n) = \sum_{x^n \in A_\epsilon^{(n)}} p(x^n) l(x^n) + \sum_{x^n \in A_\epsilon^{(n)c}} p(x^n) l(x^n) \\
 &\leq \sum_{x^n \in A_\epsilon^{(n)}} p(x^n) [n(H+\epsilon) + 2] + \sum_{x^n \in A_\epsilon^{(n)c}} p(x^n) (n \log |X| + 2) \\
 &\leq n(H+\epsilon) + \epsilon n \log |X| + 2 \\
 &= n(H + \underbrace{\epsilon \log |X| + \frac{2}{n}}_{2\epsilon'})
 \end{aligned}$$

The extra term (summed into ϵ') can be made arbitrarily small by proper choice of ϵ and n .

Theorem: Let X^n be i.i.d. $\sim p(x)$. Let $\epsilon > 0$.

Then, there is a code that maps sequences x^n of length n into binary strings such that the mapping is one-to-one, therefore, lossless and

$$E\left[\frac{1}{n} l(x^n)\right] \leq H(X) + \epsilon$$

Entropy Rate of a Stochastic Process

In order to extend the notion of entropy, which we have only defined for i.i.d. sources to sources with memory, we define the entropy rate:

Definition: The entropy rate of a random process X is defined as:

$$H(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, \dots, X_n)$$

Another definition that we show later to be equivalent to the above is given as

$$H'(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X_n | X_{n-1}, X_{n-2}, \dots, X_1)$$

before proving the equivalence of $H(\mathcal{X})$ and $H'(\mathcal{X})$ let's find the entropy rate of the memoryless processes.

For an i.i.d. process:

$$H(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{H(X_1, \dots, X_n)}{n} = \lim_{n \rightarrow \infty} \frac{n H(X_1)}{n} = H(X_1)$$

Similarly, for a source that is memoryless but not with identically distributed samples, we have:

$$H(X_1, \dots, X_n) = \sum_{i=1}^n H(X_i)$$

In this case, it is possible that $\lim_{n \rightarrow \infty} \frac{H(X_1, \dots, X_n)}{n}$ does not exist.

As an example, take a binary sequence with

$$p_i = P(X_i = 1) = \begin{cases} 0.5 & 2^k \leq \log \log i \leq 2^{k+1} \\ 0 & 2^{k+1} \leq \log \log i \leq 2^{k+2} \end{cases}$$

for $k = 0, 1, 2, \dots$

In this example $H(X_i)$ will be 1 for certain period of time followed by a longer period of time when it $H(X_i) = 0$. So, the running average oscillates between 0 and 1.

To show that $H(\mathcal{X}) = H'(\mathcal{X})$, we first prove that the limit $H'(\mathcal{X}) = \lim_{n \rightarrow \infty} (H(X_n | X_{n-1}, \dots, X_1))$ exists.

Theorem: For a stationary stochastic process, $H(X_n | X_{n-1}, \dots, X_1)$ is decreasing in n and has limit $H'(\mathcal{X})$.

Proof:

$$\begin{aligned} H(X_{n+1} | X_1, \dots, X_n) &\leq H(X_{n+1} | X_n, \dots, X_2) \\ &= H(X_n | X_{n-1}, \dots, X_1) \end{aligned}$$

So $H(X_n | X_{n-1}, \dots, X_1)$ is a decreasing sequence of non-negative numbers and, therefore, must have a limit $H'(\mathcal{X})$.

Another theorem that we need in order to prove the equivalence of $H(\mathcal{X})$ and $H'(\mathcal{X})$ is:

Theorem: If $a_n \rightarrow a$ then $b_n = \frac{1}{n} \sum_{i=1}^n a_i \rightarrow a$.

Proof: since $a_n \rightarrow a$, there is some $N(\epsilon)$ s.t., $|a_n - a| < \epsilon$, $n \geq N$

$$|b_n - a| = \left| \frac{1}{n} \sum_{i=1}^n (a_i - a) \right| \leq \frac{1}{n} \sum_{i=1}^n |a_i - a|$$

$$= \frac{1}{n} \sum_{i=1}^{N(\epsilon)} |a_i - a| + \frac{1}{n} \sum_{i=N(\epsilon)}^n |a_i - a|$$

$$\leq \frac{1}{n} \sum_{i=1}^{N(\epsilon)} |a_i - a| + \frac{n - N(\epsilon)}{n} \epsilon \quad \forall n \geq N(\epsilon)$$

as $n \rightarrow \infty$, the first term $\rightarrow 0$ so $b_n \rightarrow a$.

Now, we are ready to prove the equivalence of $H(\mathcal{X})$ and $H'(\mathcal{X})$.

Theorem: For a stationary stochastic process

$$H(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{H(X_1, \dots, X_n)}{n} = \lim_{n \rightarrow \infty} H(X_n | X_{n-1}, \dots, X_1) = H'(\mathcal{X})$$

Proof:

by the chain rule

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

So:

$$H(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, \dots, X_n) = \lim_{n \rightarrow \infty} H(X_n | X_{n-1}, \dots, X_1) = H'(\mathcal{X}).$$

Here, we have used the fact that the limit of ~~each~~ the terms in the summation is $H'(\mathcal{X})$.

The significance of the entropy rate is due to the AEP property:

$$-\frac{1}{n} \log P(X_1, \dots, X_n) \rightarrow H(\mathcal{X})$$

with probability 1.

Definition: ~~*~~(Markov chain): A random process $X_1, X_2, \dots, X_n, \dots$ is said to form a Markov chain if: