

### 2.1. Classical Information Sources : Shannon's Noiseless Theorem

Classical information sources are emitters of letters  $x_i \in \mathcal{X}$  from an alphabet  $\mathcal{X}$  consisting of  $m = \#(\mathcal{X}) = \text{card}(\mathcal{X})$  symbols.

Any string consisting of  $n$  successive letters  $x_{i(n)} = x_{i_1} x_{i_2} \dots x_{i_n}$ ,  $x_{i_j} \in \mathcal{X}$ , is a word of length  $n$

emitted by the source:  $x_{i(n)} \in \mathcal{X}^{(n)} = \underbrace{\mathcal{X} \times \mathcal{X} \times \dots \times \mathcal{X}}_{n \text{ times}}$ .  
 $i^{(n)} = i_1 i_2 \dots i_n$

There are  $\#(\mathcal{X}^{(n)}) = m^n$  such words but not all of them occur with the same probability (frequency) during a speech.

**Definition 2.1.1.** A classical information source can be described as a stochastic process  $\{X_j\}_{j=1}^{\infty}$  consisting of random variables  $X_j$  with values  $x_{i_j} \in \mathcal{X}$  occurring with joint probabilities  $P[x_{i_1} x_{i_2} \dots x_{i_n}] = \text{Prob}(X_1 = x_{i_1}, X_2 = x_{i_2}, \dots, X_n = x_{i_n})$  for all  $n \geq 1$ .

Remark: the probability distributions  $\Pi_{[1, n]} = \{ P_{[1, n]}(x_{[1, n]}) \}_{x_{[1, n]} \in \mathcal{X}^{(n)}}$

must be compatible: 
$$\begin{aligned} \sum_{i_n} P_{[1, n]}(x_{[1, n]}) &= \sum_{i_n} P_{[1, n]}(x_{i_1} x_{i_2} \dots x_{i_n}) \\ &= P_{[1, n-1]}(x_{i_1} x_{i_2} \dots x_{i_{n-1}}) \end{aligned}$$

Notice that, in general, 
$$\begin{aligned} \sum_{i_1} P_{[1, n]}(x_{[1, n]}) &= \sum_{i_1} P_{[1, n]}(x_{i_1} x_{i_2} \dots x_{i_n}) \\ &= P_{[2, n]}(x_{i_2} x_{i_3} \dots x_{i_n}) \\ &\neq P_{[1, n-1]}(x_{i_2} x_{i_3} \dots x_{i_n}) \end{aligned}$$

**Definition 2.1.2.**

If a source is described by probability distributions  $\Pi_{[1, n]} = \{ P_{[1, n]}(x_{[1, n]}) \}_{x_{[1, n]} \in \mathcal{X}^{(n)}}$

such that 
$$\sum_{i_1} P_{[1, n]}(x_{i_1} x_{i_2} \dots x_{i_n}) = P_{[1, n-1]}(x_{i_2} x_{i_3} \dots x_{i_n})$$

the source is called stationary.

Remark : stationary sources are such that the probability of any word depends only on the word itself and not on which one of the source it started to be emitted.

Example 2.1.1 Bernoulli sources

They are described by independent, identically distributed (i.i.d.)

random variables : 
$$P_{[1,n]}(x_{i_1} x_{i_2} \dots x_{i_n}) = \prod_{j=1}^n p(x_{i_j})$$

Compatibility : 
$$\sum_{i_n} P_{[1,n]}(x_{i_1} x_{i_2} \dots x_{i_n}) = \prod_{j=1}^{n-1} p(x_{i_j}) = P_{[1,n-1]}(x_{i_1} x_{i_2} \dots x_{i_{n-1}})$$

Stationarity : 
$$\sum_{i_1} P_{[1,n]}(x_{i_1} x_{i_2} \dots x_{i_n}) = \prod_{j=2}^n p(x_{i_j}) = P_{[1,n-1]}(x_{i_2} x_{i_3} \dots x_{i_n})$$

The letters from a Bernoulli source are all uncorrelated with each other.

Remark : given the joint probabilities  $P_{[1,n]}(x_{i_1} x_{i_2} \dots x_{i_n})$ ,  
one defines the conditional probabilities

$$P(x_{i_n} | x_{i_1} x_{i_2} \dots x_{i_{n-1}}) = \frac{P_{[1,n]}(x_{i_1} x_{i_2} \dots x_{i_n})}{P_{[1,n-1]}(x_{i_1} x_{i_2} \dots x_{i_{n-1}})}$$

of the last letter on the preceding ones. Then,

$$\begin{aligned} P_{[1,n]}(x_{i_1} x_{i_2} \dots x_{i_n}) &= P(x_{i_n} | x_{i_1} \dots x_{i_{n-1}}) P_{[1,n-1]}(x_{i_1} x_{i_2} \dots x_{i_{n-1}}) \\ &= P(x_{i_n} | x_{i_1} \dots x_{i_{n-1}}) P(x_{i_{n-1}} | x_{i_1} \dots x_{i_{n-2}}) \\ &\quad \times P_{[1,n-2]}(x_{i_1} \dots x_{i_{n-2}}) \\ &= \left( \prod_{j=2}^n P(x_{i_j} | x_{i_1} \dots x_{i_{j-1}}) \right) \times P_{\{1\}}(x_{i_1}) \end{aligned}$$

Example 2.1.2

Markov sources

The conditional probabilities depend only on the preceding letter:

$$P(x_{i_n} | x_{i_1} \dots x_{i_{n-1}}) = P(x_{i_n} | x_{i_{n-1}})$$

$$P_{[1,n]}(x_{i_1} x_{i_2} \dots x_{i_n}) = P(x_{i_n} | x_{i_{n-1}}) P(x_{i_{n-1}} | x_{i_{n-2}}) \dots P(x_{i_2} | x_{i_1}) P_{[1,3]}(x_{i_1})$$

Compatibility :  $\sum_{i_n} P_{[1,n]}(x_{i_1} \dots x_{i_n}) = P_{[1,n-1]}(x_{i_1} \dots x_{i_{n-1}}) \Leftrightarrow \boxed{\sum_{i_n} P(x_{i_n} | x_{i_{n-1}}) = 1}$

Stationarity :  $\sum_{i_1} P_{[1,n]}(x_{i_1} \dots x_{i_n}) = P_{[1,n-1]}(x_{i_2} \dots x_{i_n}) \Leftrightarrow \boxed{\sum_{i_1} P(x_{i_2} | x_{i_1}) P_{[1,3]}(x_{i_1}) = P_{[1,3]}(x_{i_2})}$

The transition matrix  $[T_{ij} := P(x_i | x_j)] =: T$  must be such that the sum of any of its ~~rows~~ <sup>columns</sup> is unity, because of compatibility.

Stationarity requires that  $\boxed{T |\pi\rangle = |\pi\rangle}$  where  $|\pi\rangle = \begin{pmatrix} p(x_1) \\ \vdots \\ p(x_n) \end{pmatrix}$ .

Question : what is the limit to which information can be reliably compressed, that is, compressed in a manner such that it can be recovered later with arbitrarily low probability of error ?

**Definition 2.1.3**

A source emitting words  $x_{i(n)} \in \mathcal{X}^{(n)}$  is reliably encodable at rate  $R > 0$

if, for any  $n$ , there exists a set  $A^{(n)} \subseteq \mathcal{X}^{(n)}$  such that

$$\#(A^{(n)}) \leq 2^{nR} \quad (1)$$

and

$$\begin{aligned} \lim_n \text{Prob}(A^{(n)}) &= 1 & (2) \\ \parallel \\ \lim_n \sum_{x_{i(n)} \in A^{(n)}} P(x_{i(n)}) & \end{aligned}$$

**Definition 2.1.4.**

The information <sup>rate</sup>  $\sqrt{H}$  of a given source is the smallest reliable encoding rate:  $H = \inf \{R : R \text{ reliable}\}$

Remark :  $\#(A^{(n)}) \leq 2^{nR}$  means that we can specify the elements of  $A^{(n)}$  by means of at most  $R$  bits per symbol, to be compared with the  $\log_2 m$  bits per symbol in relation to the whole set  $\mathcal{X}^{(n)}$ .

$\text{Prob}(A^{(n)}) \xrightarrow[n \rightarrow \infty]{} 1$  means that if we restrict the coding to the set  $A^{(n)}$  the probability of error, that is of non coding a relevant string in  $(A^{(n)})^c = \mathcal{X}^{(n)} \setminus A^{(n)}$  vanishes with  $n \rightarrow +\infty$  :

$$\text{Prob}((A^{(n)})^c) = 1 - \text{Prob}(A^{(n)})$$

Theorem 2.1.1 For a source with  $\#(\mathcal{X}) = m$   $0 \leq H \leq \log_2 m$

Proof:  $H \geq 0$  since  $H = \inf\{R : R \text{ reliable}\}$  and  $R > 0$ .  
 $H \leq \log_2 m$  since  $A^{(n)} \subseteq \mathcal{X}^{(n)} \Rightarrow \#(A^{(n)}) \leq 2^{n \log_2 m}$ .

Exercise 2.1.1

Show how to reach  $H=0$  and  $H=\log_2 m$ .

93

- $H=0$ : uniform source.  $P(x_{i(n)})=0$  unless

$$x_{i(n)} = x_{i(n)}^{\text{unif}} = \underbrace{x_i x_i \dots x_i}_{n \text{ times}}. \text{ Then, } \text{Prob}(A^{(n)} \subseteq \mathcal{X}^{(n)}) = 0 \text{ unless}$$

$$A^{(n)} = \{ x_{i(n)}^{\text{unif}} : x_i \in A \subseteq \mathcal{X} \} \text{ and } \boxed{\text{Prob}(A^{(n)}) = 1 \text{ iff } A = \mathcal{X}}$$

Then,  $\#(A^{(n)}) = m < 2^{nR}$  for all  $R > 0$  with  $n$  large enough

So  $\boxed{\text{inf}_f R = 0}$ .

- $H = \log_2 m$ : equiprobable Bernoulli source.

$$P(x_{i(n)}) = \prod_{j=1}^n P(x_{i_j}) = \frac{1}{m^n}$$

Suppose  $A^{(n)}$  is such that  $\#(A^{(n)}) \leq 2^{Rn}$ .

$$\text{Then, } \text{Prob}(A^{(n)}) = \sum_{i^{(n)}: x_{i(n)} \in A^{(n)}} P(x_{i(n)}) = \frac{\#(A^{(n)})}{m^n} \leq 2^{n(R - \log_2 m)} \xrightarrow{n} 0$$

whenever  $R < \log_2 m$ . Therefore,  $\boxed{H = \inf_R \{R: R \text{ reliable}\} = \log_2 m}$

**Lemma 2.1.1**

Set  $D_m(R) := \max \{ \text{Prob}(A) : \#(A) \leq 2^{nR} \}$ .

For any  $\epsilon > 0$ , the information rate  $H$  satisfies

$\lim_n D_n(H+\epsilon) = 1$  (1)

and

$H > 0 \Rightarrow \lim_n D_n(H-\epsilon) \neq 1$  (2)

Proof:

$H+\epsilon$  is a reliable encoding rate:  $\exists \{A^{(n)}\}_{n \geq 1}$  such that  $\#(A^{(n)}) \leq 2^{n(H+\epsilon)}$  and  $\text{Prob}(A^{(n)}) \xrightarrow{n} 1$ .

$1 \geq D_n(H+\epsilon) \geq \text{Prob}(A^{(n)}) \xrightarrow{n} 1 \Rightarrow \lim_n D_n(H+\epsilon) = 1$

$H > 0$ :  $\exists \epsilon > 0$ :  $0 < H-\epsilon$  is not a reliable encoding rate

Choose  $C_{H-\epsilon}^{(n)} \subseteq \mathcal{X}^{(n)}$  such that  $D_n(H-\epsilon) = \text{Prob}(C_{H-\epsilon}^{(n)})$ .

Then,  $\#(C_{H-\epsilon}^{(n)}) \leq 2^{n(H-\epsilon)}$  but  $\lim_n \text{Prob}(C_{H-\epsilon}^{(n)}) \neq 1$

**Definition 2.1.5.**

The log-likelihood per source letter of  $x_{i(n)}$  is

$$\xi_n(x_{i(n)}) = -\frac{1}{n} \log P(x_{i(n)}) \text{ if } P(x_{i(n)}) > 0$$

$$\xi_n(x_{i(n)}) = 0 \text{ if } P(x_{i(n)}) = 0$$

Remark :  $\xi_n(x_{i(n)})$  is a random variable

**Lemma 2.1.2**

For all  $R > 0$  and  $\epsilon > 0$  :  $\text{Prob}(\xi_n \leq R) \leq D_n(R) \leq \text{Prob}(\xi_n \leq R + \epsilon) + 2^{-n\epsilon}$

Proof: Set  $B_R^{(n)} = \{x_{i(n)} \in \mathcal{X}^{(n)} : P(x_{i(n)}) \geq 2^{-nR}\}$

$$= \{x_{i(n)} \in \mathcal{X}^{(n)} : \xi_n(x_{i(n)}) = -\frac{1}{n} \log P(x_{i(n)}) \leq R\}$$

• Left inequality

$$1 \geq \text{Prob}(B_R^{(n)}) = \sum_{i(n): x_{i(n)} \in B_R^{(n)}} P(x_{i(n)}) \geq \#(B_R^{(n)}) 2^{-nR} \Rightarrow \#(B_R^{(n)}) \leq 2^{nR}$$

Then,  $D_n(R) \geq \text{Prob}(B_R^{(n)})$  since  $D_n(R) = \max \{ \text{Prob}(A) : \#(A) \leq 2^{nR} \}$

• Right inequality: let  $C_R^{(n)}$  be such that  $D_n(R) = \text{Prob}(C_R^{(n)})$

$$\begin{aligned}
D_n(R) &= \text{Prob}(C_R^{(n)}) = \text{Prob} \left\{ X_{i^{(n)}} \in C_R^{(n)} : \sum_n (X_{i^{(n)}}) \leq R + \epsilon \right\} \\
&\quad + \text{Prob} \left\{ X_{i^{(n)}} \in C_R^{(n)} : \sum_n (X_{i^{(n)}}) > R + \epsilon \right\} \\
&\leq \text{Prob} \left( \sum_n \leq R + \epsilon \right) + \sum_{i^{(n)}: X_{i^{(n)}} \in C_R^{(n)}, P(X_{i^{(n)}}) < 2^{-n(R+\epsilon)}} P(X_{i^{(n)}}) \\
&\leq \text{Prob} \left( \sum_n \leq R + \epsilon \right) + 2^{-n(R+\epsilon)} \#(C_R^{(n)}) \\
&\leq \text{Prob} \left( \sum_n \leq R + \epsilon \right) + 2^{-n(R+\epsilon)} 2^{nR}
\end{aligned}$$

**Definition 2.1.6**

A sequence of random variables  $X_n$  converges in probability to a constant  $x$  ( $X_n \xrightarrow{P} x$ )

if for any  $\epsilon > 0$   $\lim_n \text{Prob}(|X_n - x| \geq \epsilon) = 0$

Example 2.1.3

Asymptotic Equipartition Property

Suppose  $\sum_n \frac{P}{n} \rightarrow H$ , then  $\forall \eta > 0, \epsilon > 0$  there exists  $n_0(\epsilon, \eta)$

such that  $n > n_0(\epsilon, \eta)$  implies  $\text{Prob}(|\xi_n - H| \geq \epsilon) < \eta$

Then,  $1 - \eta \leq \text{Prob}(|\xi_n - H| < \epsilon) = \text{Prob}\left\{x_{i(n)} : H - \epsilon < \sum_n(x_{i(n)}) < H + \epsilon\right\}$

$$1 - \eta \leq \text{Prob}\left\{x_{i(n)} : 2^{-n(H+\epsilon)} < p(x_{i(n)}) < 2^{-n(H-\epsilon)}\right\}$$

call  $T_{\epsilon, \eta}^{(n)} := \left\{x_{i(n)} : 2^{-n(H+\epsilon)} < p(x_{i(n)}) < 2^{-n(H-\epsilon)}\right\}$

- $\text{Prob}(T_{\epsilon, \eta}^{(n)}) \geq 1 - \eta$

- $1 \geq \text{Prob}(T_{\epsilon, \eta}^{(n)}) \geq \#(T_{\epsilon, \eta}^{(n)}) 2^{-n(H+\epsilon)} \Rightarrow \#(T_{\epsilon, \eta}^{(n)}) \leq 2^{n(H+\epsilon)}$

- $1 - \eta \leq \text{Prob}(T_{\epsilon, \eta}^{(n)}) \leq \#(T_{\epsilon, \eta}^{(n)}) 2^{-n(H-\epsilon)} \Rightarrow \#(T_{\epsilon, \eta}^{(n)}) \geq (1 - \eta) 2^{n(H-\epsilon)}$

Remark :

$T_{\epsilon, \eta}^{(n)}$  is called a typical subset and the strings it contains typical strings.  $T_{\epsilon, \eta}^{(n)}$  is a high probability subset and one can decide to encode typical strings by means of roughly  $H$  bits per symbol.

Untypical strings in  $\mathcal{X}^{(n)} \setminus T_{\epsilon, \eta}^{(n)} =: (T_{\epsilon, \eta}^{(n)})^c$  can be all encoded by a same string, which would then correspond to errors. However, the probability of such an error is bounded by  $\text{Prob}((T_{\epsilon, \eta}^{(n)})^c) \leq \eta$

Theorem 2.1.2

Shannon's Noiseless Coding Theorem

If  $\bar{X}_n(x_i^{(n)}) = -\frac{1}{n} \log P(x_i^{(n)})$  converges to a constant  $\xi$ , then  $\xi = H$ , the information rate of the source.

Proof: If  $\sum_m P \rightarrow \xi$  then  $\xi \geq 0$  since  $\sum_m (X_{\epsilon(m)}) \geq 0$ .

(33)

From Lemma 2.1.2:  $\forall \epsilon > 0$

$$D_n(\xi + \epsilon) \geq \text{Prob}(\sum_m \leq \xi + \epsilon) \geq \text{Prob}(\xi - \epsilon \leq \sum_m \leq \xi + \epsilon)$$

||

$$\text{Prob}(|\sum_m - \xi| \leq \epsilon)$$

||

$$\boxed{H \leq \xi} \quad \leftarrow \quad 1 - \frac{1}{n} \leftarrow 1 - \text{Prob}(|\sum_m - \xi| > \epsilon)$$

Notice that  $\boxed{\xi = 0 \Rightarrow H = 0}$ .

If  $\xi > 0$ , for any  $\epsilon > 0$ , Lemma 2.1.2 yields

$$D_n(\xi - \epsilon) \leq \text{Prob}\left(\sum_m \leq \left(\xi - \epsilon\right) + \frac{\epsilon}{2}\right) + 2^{-n\epsilon/2}$$

$$\leftarrow \text{Prob}\left(\xi - \sum_m \geq \frac{\epsilon}{2}\right) + 2^{-n\epsilon/2}$$

$$\leq \text{Prob}\left(|\xi - \sum_m| \geq \frac{\epsilon}{2}\right) + 2^{-n\epsilon/2} \xrightarrow{n} 0$$

Then, Lemma 2.1.1 yields  $\boxed{H \geq \xi}$